

Rank adaptive tensor recovery based model reduction for partial differential equations with high-dimensional random inputs

Kejun Tang^a, Qifeng Liao^{b,*}

^a*School of Information Science and Technology, ShanghaiTech University, Shanghai, China*

^b*School of Information Science and Technology, ShanghaiTech University, Shanghai, China*

Abstract

This work proposes a systematic model reduction approach based on rank adaptive tensor recovery for partial differential equation (PDE) models with high-dimensional random parameters. Since the standard outputs of interest of these models are discrete solutions on given physical grids which are high-dimensional, we use kernel principal component analysis to construct stochastic collocation approximations in reduced dimensional spaces of the outputs. To address the issue of high-dimensional random inputs, we develop a new efficient rank adaptive tensor recovery approach to compute the collocation coefficients. Novel efficient initialization strategies for non-convex optimization problems involved in tensor recovery are also developed in this work. We present a general mathematical framework of our overall model reduction approach, analyze its stability, and demonstrate its efficiency with numerical experiments.

Keywords: tensor recovery; model reduction; PDEs; uncertainty quantification

1. Introduction

During the last few decades there has been a rapid development in surrogate and reduced order modelling for PDE systems with random inputs. The PDE systems are fundamental mathematical models describing complex physical and engineering problems, which can involve multiple disciplines, a large number of input parameters, and multiple sources of uncertainty. A main challenge of surrogate modelling for these PDE models is the so-called curse of dimensionality. First, due to the high complexity of practical problems, the random input parameters are typically high-dimensional. Second, the standard output of these PDE models is the spatial fields (e.g., temperature, pressure and velocity), and their fine resolution representation requires a large number of degrees of freedom, which make the output high-dimensional.

A type of widely used surrogate modelling approach for these PDE models is the stochastic spectral methods [1, 2, 3, 4, 5], while the high dimensionality of the random inputs causes difficulties in applying them. To alleviate the difficulty, modifications of these methods have been actively introduced by exploiting certain properties of the underlying problem. For example, sparse (generalized) polynomial chaos (gPC) expansions [6, 7, 8, 9, 10, 11] are developed through using the sparsity in spectral approximations. Moreover, the stochastic collocation method [3, 12, 4] is reformulated as a tensor style quadrature problem in [13], which shows that the corresponding collocation

*Corresponding author

Email addresses: tangkj@shanghaitech.edu.cn (Kejun Tang), liaofq@shanghaitech.edu.cn (Qifeng Liao)

coefficients can be efficiently computed through tensor recovery techniques. On the other hand, the high dimensionality in outputs poses challenges in both surrogate modelling and data storage. The surrogates proposed to resolve high-dimensional inputs as discussed above are typically restricted to problems with a single output. Naively extending them to high-dimensional outputs (building independent surrogates for multiple outputs) is computationally infeasible. For making progress, dimension reduction methods for the outputs gain a lot of interests. For example, principal component analysis (PCA) and kernel component analysis (kPCA) methods are successfully established for Gaussian process surrogates [14, 15]. Especially, since kPCA captures highly nonlinear low-rank structures in the output space, it can provide dramatically tight representation of the outputs [16, 15].

In this work, we focus on tensor recovery based stochastic collocation. As discussed in [13], gPC coefficients in stochastic collocation can be computed through inner products of weight tensors and data tensors (see section 2.2 for details), where the weight tensors are given but the data tensors are expensive to obtain. Instead of directly evaluating the expensive data tensor, tensor recovery here is to use a small number of entries of the data tensor to recover the whole tensor [17, 18, 19]. A popular recovery strategy is developed in [17] based on canonical polyadic (CP) decomposition. In this recovery approach, the CP rank of the underlying tensor needs to be known a priori, which limits its application to our PDE models where the corresponding CP ranks are not given. For this purpose, we develop a novel rank adaptive tensor recovery (RATR) approach, which do not require any prior information for the tensor ranks. Moreover, as this kind of tensor recovery procedure requires solving a non-convex optimization problem [13], initialization strategies for this kind of optimization problem are crucial for successful recovery. In our RATR approach, new efficient initialization strategies are proposed based on a hierarchical rank-one updating procedure, and their stability is theoretically proven in this work.

The aim of this paper is to develop a systematic model reduction framework to curb this challenging high-dimensional input-output problem, and our overall procedure is as follows. First, kPCA is conducted for the outputs, which gives their reduced-dimensional representations. After that, for each kPCA mode, RATR based stochastic collocation is proposed to construct sparse gPC expansions for each kPCA mode. The inverse mapping method introduced in [16, 15] is finally adopted to construct an overall estimates of the outputs in the high-dimensional space. To summarize, the main contributions of this work are three-fold: first, stochastic collocation methods are reformulated with manifold learning for high-dimensional outputs; second, a novel rank adaptive tensor recovery (RATR) approach is proposed to recover tensors without knowing their ranks a priori; third, new efficient initialization strategies for RATR are proposed and their stability is analyzed.

The rest of the paper is organized as follows. In the next section, we formulate stochastic collocation methods for stochastic PDEs based on manifold learning. Details of tensors and standard tensor recovery approaches are introduced in section 3. Our main algorithms and analysis for RATR and the overall RATR-collocation surrogate are presented in section 4. In section 5, we demonstrate the efficiency of our RATR-collocation approach for stochastic diffusion and incompressible flow problems. Finally section 6 concludes the paper.

2. Problem setting and stochastic collocation based on manifold learning

Let D denote a spatial domain (in \mathbb{R}^2 or \mathbb{R}^3) which is bounded, connected and with a polygonal boundary ∂D , and x denote a spatial variable. Let ξ be a vector which collects a finite number of random variables. The dimension of ξ is denoted by d , i.e., we write $\xi = [\xi_1, \dots, \xi_d]^T$. The probability density function of ξ is denoted by $\pi(\xi)$. In this paper, we restrict our attention to the situation that ξ has a bounded and connected support. Without loss of generality, we next assume the support of ξ to be I^d where $I := [-1, 1]$, since any bounded connected domain in \mathbb{R}^d can be mapped to I^d . The physics of problems considered in this paper are governed by a PDE over the spatial domain D and boundary conditions on the boundary ∂D . This PDE problem is stated as: find $u(x, \xi) : D \times I^d \rightarrow \mathbb{R}$, such that

$$\mathcal{L}(x, \xi; u(x, \xi)) = f(x) \quad \forall (x, \xi) \in D \times I^d, \quad (1)$$

$$\mathfrak{b}(x, \xi; u(x, \xi)) = g(x) \quad \forall (x, \xi) \in \partial D \times I^d, \quad (2)$$

where \mathcal{L} is a partial differential operator and \mathfrak{b} is a boundary operator, both of which can have random coefficients. f is the source function and g specifies the boundary conditions. We also define output quantities of interest. For each realization of ξ , if the deterministic version of (1)–(2) is solved using a high-fidelity numerical scheme (simulator), for example finite element and difference methods, a natural definition of the output is the discrete solution. A high-fidelity discrete solution is also called a *snapshot* and can be represented as $\mathbf{y} = [u(x^{(1)}, \xi), \dots, u(x^{(N_h)}, \xi)]^T \in \mathbb{R}^{N_h}$, where $u(x^{(i)}, \xi)$, $i = 1, \dots, N_h$ denotes the value of $u(x, \xi)$ at a specified location on a spatial grid and N_h refers to the spatial degrees of freedom. The manifold consisting of all snapshots is denoted by $\mathcal{M} \subset \mathbb{R}^{N_h}$, and it is assumed to be smooth. A PDE simulator can be viewed as a mapping $\chi : I^d \rightarrow \mathcal{M}$, where I^d and $\mathcal{M} \in \mathbb{R}^{N_h}$ are the input space and the output manifold respectively, and we denote it as $\chi(\xi) = \mathbf{y} = [u(x^{(1)}, \xi), \dots, u(x^{(N_h)}, \xi)]^T \in \mathcal{M}$ for an arbitrary realization of the input $\xi \in I^d$.

The goal of this study is to build surrogates for conducting uncertainty qualification (UQ) of the output \mathbf{y} , given limited training data points $\mathbf{y}^{(j)} = \chi(\xi^{(j)})$ for $j = 1 : N_t$ where N_t is the size of a training data set. We focus on the challenging situations that the input and the output are both high-dimensional. To make progress, we reformulate the stochastic collocation surrogates [4, 3] based on manifold learning and tensor recovery quadrature. Manifold learning gives a reduced dimension representation for the output space through kernel principal component analysis (kPCA) and inverse mappings [20, 21, 15], and tensor recovery provides estimates of collocation coefficients associated with high-dimensional random parameters through exploiting low rank structures in these coefficients [13]. The rest of this section is to discuss the manifold learning based collocation and the setting of tensor formulation, while detailed tensor recovery methods and our new rank adaptive schemes are presented in the next two sections.

2.1. Kernel principal component analysis (kPCA)

To simplify the presentation, the given training data are denoted by $\mathbf{y}^{(j)} = \chi(\xi^{(j)})$ for $j = 1, \dots, N_t$. Following [20, 22, 15], the kernel principal component analysis (kPCA) proceeds through two steps: mapping the training data to a higher-dimensional feature space, and performing linear principal component analysis (PCA) in the feature space.

Denoting the feature space by \mathcal{F} , we define a mapping $\Gamma : \mathcal{M} \rightarrow \mathcal{F}$, which maps each training data point $\mathbf{y}^{(j)} \in \mathcal{M}$ to $\Gamma(\mathbf{y}^{(j)}) \in \mathcal{F}$ for $j = 1, \dots, N_t$. A covariance matrix of the mapped data is defined as

$$\mathbf{C}_{\mathcal{F}} := \frac{1}{N_t} \sum_{j=1}^{N_t} \tilde{\Gamma}(\mathbf{y}^{(j)}) \tilde{\Gamma}(\mathbf{y}^{(j)})^T, \quad (3)$$

where $\tilde{\Gamma}(\mathbf{y}^{(j)}) = \Gamma(\mathbf{y}^{(j)}) - \bar{\Gamma}$ and $\bar{\Gamma} = (1/N_t) \sum_{j=1}^{N_t} \Gamma(\mathbf{y}^{(j)})$. Eigenvectors of $\mathbf{C}_{\mathcal{F}}$ can give a new basis to represent the mapped data, and the eigenvectors associated with dominate eigenvalues can provide an effective reduced dimensional representation for them.

However, the mapping Γ in practice is typically defined implicitly through kernel functions, and the eigenvectors of $\mathbf{C}_{\mathcal{F}}$ are always replaced by eigenvectors of some centred kernel matrices. A kernel function in this setting is a mapping from $\mathbb{R}^{N_h} \times \mathbb{R}^{N_h}$ to \mathbb{R} , which is denoted by $\mathbf{k}(\cdot, \cdot)$. The kernel matrix associated with $\mathbf{k}(\cdot, \cdot)$ is denoted by $\mathbf{K} \in \mathbb{R}^{N_h \times N_h}$, of which each entry is defined as

$$\mathbf{K}_{ij} = \mathbf{k}(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) \quad \text{for } i, j = 1, 2, \dots, N_t.$$

A standard choice of the kernel function is the Gaussian kernel

$$\mathbf{k}(\mathbf{y}^{(i)}, \mathbf{y}^{(j)}) = \exp\left(-\frac{\|\mathbf{y}^{(i)} - \mathbf{y}^{(j)}\|_2^2}{2\sigma_g^2}\right),$$

where σ_g is the bandwidth parameter. The centred kernel matrix is defined as

$$\tilde{\mathbf{K}} := \left(\mathbf{K} - \mathbf{1}_{\frac{1}{N_t}} \mathbf{K} - \mathbf{K} \mathbf{1}_{\frac{1}{N_t}} + \mathbf{1}_{\frac{1}{N_t}} \mathbf{K} \mathbf{1}_{\frac{1}{N_t}}\right),$$

where $\mathbf{1}_{\frac{1}{N_t}}$ denotes the matrix with all entries equaling to $1/N_t$.

Let $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_{N_t}] \in \mathbb{R}^{N_t \times N_t}$ collect the eigenvectors of $\tilde{\mathbf{K}}$ associated with eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_{N_t}$. Basis functions of the mapped data in \mathcal{F} are defined as $\omega_e := \sum_{j=1}^{N_t} \tilde{\alpha}_{je} \tilde{\Gamma}(\mathbf{y}^{(j)})$, where $\tilde{\alpha}_{je} = \alpha_{je} / \sqrt{\lambda_e}$ and α_{je} is the j -th component of $\boldsymbol{\alpha}_e$, for $e = 1, \dots, N_t$. A mapped training point $\tilde{\Gamma}(\mathbf{y}^{(j)})$ for $j = 1, \dots, N_t$ can be represented as $\tilde{\Gamma}(\mathbf{y}^{(j)}) = \sum_{e=1}^{N_t} \tilde{\gamma}_e(\mathbf{y}^{(j)}) \omega_e$, where each coefficient $\tilde{\gamma}_e(\mathbf{y}^{(j)})$ is computed through

$$\tilde{\gamma}_e(\mathbf{y}^{(j)}) = \sum_{i=1}^{N_t} \tilde{\alpha}_{ie} \tilde{\mathbf{K}}_{ij}. \quad (4)$$

To result in dimension reduction, the first N_r dominant eigenvectors of $\tilde{\mathbf{K}}$ are selected as the principal components, with the criterion $(\sum_{e=1}^{N_r} \lambda_e) / (\sum_{e=1}^{N_t} \lambda_e) > \text{tol}_{\text{PCA}}$ where tol_{PCA} is a given tolerance. The basis functions associate with the principal components are then ω_e , $e = 1, \dots, N_r$, and each mapped training data point can be approximated as $\tilde{\Gamma}(\mathbf{y}^{(j)}) \approx \sum_{e=1}^{N_r} \tilde{\gamma}_e(\mathbf{y}^{(j)}) \omega_e$.

It can be seen that the overall procedure of kPCA defines a mapping from the output manifold \mathcal{M} to the reduced feature space $\text{span}\{\omega_1, \dots, \omega_{N_r}\}$. We denote this mapping as $\kappa(\mathbf{y}) = \sum_{e=1}^{N_r} \tilde{\gamma}_e(\mathbf{y}) \omega_e$, where each coefficient $\tilde{\gamma}_e(\mathbf{y})$ is obtained from (4). The basis $\{\omega_e\}_{e=1}^{N_r}$ discussed above depends on the mapping $\tilde{\Gamma}$ which are defined implicitly. Collecting these coefficients, a reduced output vector is denoted by $\tilde{\boldsymbol{\gamma}}(\mathbf{y}) := [\tilde{\gamma}_1(\mathbf{y}), \dots, \tilde{\gamma}_{N_r}(\mathbf{y})]^T \in \mathbb{R}^{N_r}$ for any $\mathbf{y} \in \mathcal{M}$. The manifold consisting of all reduced output vectors is denoted by \mathcal{M}_r . We next denote $\boldsymbol{\gamma}(\boldsymbol{\xi}) := \tilde{\boldsymbol{\gamma}}(\mathbf{y}) \in \mathcal{M}_r \subset \mathbb{R}^{N_r}$. In summary, each training data point $\mathbf{y}^{(j)} = \chi(\boldsymbol{\xi}^{(j)}) \in \mathcal{M}$ is mapped to $\boldsymbol{\gamma}(\boldsymbol{\xi}^{(j)}) := [\gamma_1(\boldsymbol{\xi}^{(j)}), \dots, \gamma_{N_r}(\boldsymbol{\xi}^{(j)})]^T \in \mathcal{M}_r$ for $j = 1, \dots, N_t$. We next construct stochastic collocation surrogates for each component of $\boldsymbol{\gamma}(\boldsymbol{\xi})$.

2.2. Stochastic collocation

For each $\gamma_e(\xi)$, $e = 1, \dots, N_r$, a truncated generalized polynomial chaos (gPC) approximation [1, 2] can be written as

$$\gamma_e(\xi) \approx \gamma_e^{\text{gPC}}(\xi) := \sum_{\|\mathbf{i}\|_1=0}^p c_{ei} \Phi_i(\xi), \quad (5)$$

where $\mathbf{i} = [i_1, i_2, \dots, i_d]^T \in \mathbb{N}^d$ is a multi-index, $\|\mathbf{i}\|_1 = i_1 + i_2 + \dots + i_d$, and p is a given order for truncation. Denoting the set of the multi-indices by $\Upsilon := \{\mathbf{i} \in \mathbb{N}^d \text{ and } \|\mathbf{i}\|_1 = 0, \dots, p\}$, which implies that the number of basis function is $|\Upsilon| = (p+d)!/(p!d!)$. The basis functions $\{\Phi_i(\xi) | \mathbf{i} \in \Upsilon\}$ are orthogonal polynomials with respect to the density function $\pi(\xi)$

$$\langle \Phi_i(\xi), \Phi_{i'}(\xi) \rangle = \int_{I^d} \Phi_i(\xi) \Phi_{i'}(\xi) \pi(\xi) d\xi = \delta_{i,i'},$$

where δ denotes the Kronecker delta function, i.e., $\delta_{i,i'} = 1$ if \mathbf{i} is the same as \mathbf{i}' and $\delta_{i,i'} = 0$ otherwise. Each basis function $\Phi_i(\xi)$ can be expressed as the product of a set of univariate orthogonal polynomials, $\Phi_i(\xi) = \prod_{k=1}^d \phi_{i_k}(\xi_k)$, with each univariate orthogonal polynomial defined through a three term recurrence [23],

$$\begin{aligned} \phi_{j+1}(\xi) &= (\xi - \zeta_j) \phi_j(\xi) - \tau_j \phi_{j-1}(\xi), \quad j = 1, 2, \dots, p-1, \\ \phi_0(\xi) &= 0, \quad \phi_1(\xi) = 1, \end{aligned}$$

where $\zeta_j = \int_{-1}^1 \xi \pi_\xi(\xi) \phi_j^2(\xi) d\xi / \int_{-1}^1 \pi_\xi(\xi) \phi_j^2(\xi) d\xi$, $\tau_j = \int_{-1}^1 \pi_\xi(\xi) \phi_j^2(\xi) d\xi / \int_{-1}^1 \pi_\xi(\xi) \phi_{j-1}^2(\xi) d\xi$, and $\pi_\xi(\xi)$ is the marginal density function of ξ (ξ denotes a component of ξ).

According to orthogonality of the gPC basis functions, the coefficients in (5) can be computed through

$$c_{ei} = \int_{I^d} \gamma_e(\xi) \Phi_i(\xi) \pi(\xi) d\xi. \quad (6)$$

This integral can be computed through quadrature rules, and following [13] we focus on the tensor style quadrature. Let $\{\xi^{(j)}, w^{(j)}\}_{j=1}^n$ denote n quadrature nodes and weights on the interval $[-1, 1]$. The quadrature form of (6) is

$$c_{ei} = \sum_{1 \leq j_1, \dots, j_d \leq n} \gamma_e(\xi_{j_1 \dots j_d}) \Phi_i(\xi_{j_1 \dots j_d}) w_{j_1 \dots j_d}, \quad (7)$$

where

$$\xi_{j_1 \dots j_d} = [\xi_1^{(j_1)}, \xi_2^{(j_2)}, \dots, \xi_d^{(j_d)}]^T, \quad (8)$$

$$w_{j_1 \dots j_d} = w^{(j_1)} w^{(j_2)} \dots w^{(j_d)}, \quad (9)$$

for $1 \leq j_1, \dots, j_d \leq n$ are the nodes and the weights spanned by the tensor product of the one-dimensional quadrature rule.

Following [13], the quadrature form (7) can be formulated as a tensor inner product as follows. For each $e = 1, \dots, N_r$, the values $\gamma_e(\xi_{j_1 \dots j_d})$ for $1 \leq j_1, \dots, j_d \leq n$ form a d -th order *data tensor* $\mathbf{X}_e \in \mathbb{R}^{n \times \dots \times n}$, of which each entry is

$$\mathbf{X}_e(j_1, \dots, j_d) = \gamma_e(\xi_{j_1 \dots j_d}). \quad (10)$$

For each multi index \mathbf{i} with $\|\mathbf{i}\|_1 \leq p$, the values $\Phi_{\mathbf{i}}(\xi_{j_1 \dots j_d}) w_{j_1 \dots j_d}$ for $1 \leq j_1, \dots, j_d \leq n$ form a d -th order *weight tensor* $\mathbf{W}_{\mathbf{i}} \in \mathbb{R}^{n \times \dots \times n}$ with

$$\mathbf{W}_{\mathbf{i}}(j_1, \dots, j_d) = \Phi_{\mathbf{i}}(\xi_{j_1 \dots j_d}) w_{j_1 \dots j_d}. \quad (11)$$

Defining

$$\hat{\mathbf{w}}_k^{(i_k)} = [\phi_{i_k}(\xi^{(1)}) w^{(1)}, \phi_{i_k}(\xi^{(2)}) w^{(2)}, \dots, \phi_{i_k}(\xi^{(n)}) w^{(n)}]^T \in \mathbb{R}^n, \quad \text{for } k = 1, \dots, d, \quad (12)$$

each entry of $\mathbf{W}_{\mathbf{i}}$ can be written as

$$\mathbf{W}_{\mathbf{i}}(j_1, j_2, \dots, j_d) = \hat{\mathbf{w}}_1^{(i_1)}(j_1) \hat{\mathbf{w}}_2^{(i_2)}(j_2) \cdots \hat{\mathbf{w}}_d^{(i_d)}(j_d) \quad \text{for all } 1 \leq j_k \leq n, k = 1, \dots, d,$$

and $\mathbf{W}_{\mathbf{i}}$ can be expressed as

$$\mathbf{W}_{\mathbf{i}} = \hat{\mathbf{w}}_1^{(i_1)} \circ \hat{\mathbf{w}}_2^{(i_2)} \circ \cdots \circ \hat{\mathbf{w}}_d^{(i_d)} \quad \text{with } \mathbf{i} = [i_1, \dots, i_d]^T, \quad (13)$$

where “ \circ ” is the vector outer product. With the notation above, the coefficient c_{ei} in (7) can be rewritten as the tensor inner product

$$c_{ei} = \langle \mathbf{X}_e, \mathbf{W}_{\mathbf{i}} \rangle, \quad (14)$$

where the tensor inner product [24, 25] is defined as,

$$\langle \mathbf{X}_e, \mathbf{W}_{\mathbf{i}} \rangle = \sum_{j_1}^n \sum_{j_2}^n \cdots \sum_{j_d}^n \mathbf{X}_e(j_1, j_2, \dots, j_d) \mathbf{W}_{\mathbf{i}}(j_1, j_2, \dots, j_d). \quad (15)$$

The tensor norm induced by this inner product is denoted by $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$. Details of tensor decomposition and recovery are discussed in section 3 and section 4.

2.3. Inverse mapping

After the gPC approximation (5) for each $\gamma_e, e = 1, \dots, N_r$, is constructed through the above collocation procedure, the reduced output $\gamma(\xi) = \tilde{\gamma}(\mathbf{y}) = \tilde{\gamma}(\chi(\xi)) \in \mathcal{M}_r$ for an arbitrary realization of ξ can be cheaply estimated through this gPC surrogate. However, our goal is to quantify the uncertainties in the output $\mathbf{y} = \chi(\xi) \in \mathcal{M}$, which requires an inverse mapping κ^{-1} from the reduced output manifold \mathcal{M}_r to the original output manifold \mathcal{M} . Following [21, 15], an inverse mapping can be obtained through an interpolation of neighbouring points in the training data set $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(N_t)}\}$. That is, the Euclid distance between an arbitrary output $\mathbf{y} \in \mathcal{M}$ and each training point $\mathbf{y}^{(j)}$ (for $j = 1, \dots, N_t$) is first computed through

$$d_j = \sqrt{-2\sigma_g^2 \log(1 - 0.5\hat{d}_j^2)},$$

where $\hat{d}_j = 1 + \tilde{\gamma}(\mathbf{y})^T \mathbf{K} \tilde{\gamma}(\mathbf{y}) - 2\tilde{\gamma}(\mathbf{y})^T \mathbf{k}_{\mathbf{y}^{(j)}}$ are computed through the kernel function, $\mathbf{k}_{\mathbf{y}^{(j)}} = [\mathbf{k}(\mathbf{y}^{(j)}, \mathbf{y}^{(1)}), \dots, \mathbf{k}(\mathbf{y}^{(j)}, \mathbf{y}^{(N_t)})]^T$ and $\mathbf{k}(\cdot, \cdot)$ is the given kernel function. The distances $\{d_1, \dots, d_{N_t}\}$ are sorted next. Given a positive integer N_n , the indices with the smallest N_n distances are collected in a set $\mathcal{J} \subset \{1, \dots, N_t\}$, i.e., $d_j \leq d_i$ for any $j, i = 1, \dots, N_t$ with $j \in \mathcal{J}$ but $i \notin \mathcal{J}$. After that, \mathbf{y} can be approximated as

$$\mathbf{y} \approx \sum_{j \in \mathcal{J}} \frac{d_j^{-1}}{\sum_{j \in \mathcal{J}} d_j^{-1}} \mathbf{y}^{(j)}. \quad (16)$$

3. Tensor recovery based quadrature

It is clear that the main computational cost of the above collocation procedure based on manifold learning comes from generating the gPC expansion (5) for each kPCA mode $e = 1, \dots, N_r$, where evaluating each collocation coefficient requires computing a tensor inner product (14). When the input parameter ξ is high-dimensional (d is large), each data tensor $\mathbf{X}_e \in \mathbb{R}^{n \times \dots \times n}$ is large (with n^d entries). Evaluating each entry of \mathbf{X}_e requires computing a snapshot (see (10)), and it is therefore expensive to form these data tensors through computing snapshots for all entries. As an alternative, tensor recovery methods provide efficient estimates of tensors using a small number of exact entries. For forward UQ problems with a single output, when tensor ranks are given, a tensor recovery based collocation approach is developed in [13], which can be applied to construct the gPC approximation for each kPCA component (5). We here review this tensor recovery based collocation approach and provide new detailed computational cost assessments. Since computation procedures for generating the gPC surrogates for each $\gamma_e(\xi)$, $e = 1, \dots, N_r$, are identical, we generically denote the data tensor \mathbf{X}_e defined in (10) as $\mathbf{X}_{\text{exact}}$ in this section (i.e., the subscript e is temporally ignored).

3.1. Canonical polyadic (CP) decomposition

Following the presentation in [25], the CP decomposition is reviewed as follows. For a d -th order tensor $\mathbf{X} \in \mathbb{R}^{n \times \dots \times n}$, its CP decomposition is expressed as

$$\mathbf{X} = \sum_{r=1}^R \mathbf{v}_1^{(r)} \circ \mathbf{v}_2^{(r)} \circ \dots \circ \mathbf{v}_d^{(r)} \quad (17)$$

where $\mathbf{v}_k^{(r)} \in \mathbb{R}^n$ for $k = 1, \dots, d$, R is the CP rank of \mathbf{X} , and “ \circ ” is the vector outer product. The CP rank is defined as

$$R := \text{rank}(\mathbf{X}) := \min \left\{ R' \mid \mathbf{X} = \sum_{r=1}^{R'} \mathbf{v}_1^{(r)} \circ \mathbf{v}_2^{(r)} \circ \dots \circ \mathbf{v}_d^{(r)} \right\}.$$

Figure 1 shows a third-order tensor with its CP decomposition. For each $r = 1, \dots, R$, $\mathbf{v}_1^{(r)} \circ \mathbf{v}_2^{(r)} \circ \dots \circ \mathbf{v}_d^{(r)}$ in (17) is

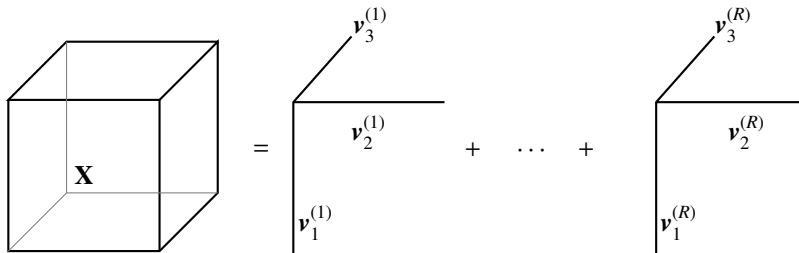


Figure 1: CP decomposition of a third-order tensor with rank R .

called a rank-one component. For each $k = 1, \dots, d$, the matrix $\mathbf{A}_k = [\mathbf{v}_k^{(1)}, \mathbf{v}_k^{(2)}, \dots, \mathbf{v}_k^{(R)}] \in \mathbb{R}^{n \times R}$ is called the k th-order factor matrix. With these factor matrices, the CP decomposition (17) can be rewritten as

$$\mathbf{X} = [[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]]. \quad (18)$$

From (13), it is clear that each weight tensor is a rank-one tensor, and we here generically denote it as $\mathbf{W} := \mathbf{w}_1 \circ \mathbf{w}_2 \circ \dots \circ \mathbf{w}_d \in \mathbb{R}^{n \times \dots \times n}$. Following [25], the inner product of \mathbf{X} , \mathbf{W} (for computing (14)) can be efficiently computed as

$$\begin{aligned}
\langle \mathbf{X}, \mathbf{W} \rangle &= \langle [[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]], \mathbf{w}_1 \circ \mathbf{w}_2 \circ \dots \circ \mathbf{w}_d \rangle \\
&= \sum_{r=1}^R \left\langle \mathbf{v}_1^{(r)} \circ \mathbf{v}_2^{(r)} \circ \dots \circ \mathbf{v}_d^{(r)}, \mathbf{w}_1 \circ \mathbf{w}_2 \circ \dots \circ \mathbf{w}_d \right\rangle \\
&= \sum_{r=1}^R \sum_{j_1, j_2, \dots, j_d} \mathbf{v}_1^{(r)}(j_1) \mathbf{v}_2^{(r)}(j_2) \dots \mathbf{v}_d^{(r)}(j_d) \mathbf{w}_1(j_1) \mathbf{w}_2(j_2) \dots \mathbf{w}_d(j_d) \\
&= \sum_{r=1}^R \left(\left(\sum_{j_1} \mathbf{v}_1^{(r)}(j_1) \mathbf{w}_1(j_1) \right) \left(\sum_{j_2} \mathbf{v}_2^{(r)}(j_2) \mathbf{w}_2(j_2) \right) \dots \left(\sum_{j_d} \mathbf{v}_d^{(r)}(j_d) \mathbf{w}_d(j_d) \right) \right) \\
&= \sum_{r=1}^R \left(\prod_{k=1}^d \langle \mathbf{v}_k^{(r)}, \mathbf{w}_k \rangle \right). \tag{19}
\end{aligned}$$

The cost for computing the tensor inner product $\langle \mathbf{X}, \mathbf{W} \rangle$ through (19) is $O(dnR)$, while the cost is $O(n^d)$ if the inner product is computed directly through its definition (15).

3.2. Missing data tensor recovery

Denoting the full index set for a d -th order tensor in $\mathbb{R}^{n \times \dots \times n}$ as $\Theta_{\text{full}} := \{[j_1, j_2, \dots, j_d]^T \in \mathbb{N}^d \mid j_k = 1, \dots, n \text{ for } k = 1, \dots, d\}$, $\Theta \subset \Theta_{\text{full}}$ denotes an observation index set, which consists of $|\Theta| \ll d^n$ indices uniformly sampled from Θ_{full} . To numbering the elements in Θ , the following sort operator is introduced.

Definition 1 (Sort operator). *For a given finite set $\Theta \subset \mathbb{N}^d$, we first sort its elements in alphabetical order: for any two different indices $\hat{\mathbf{j}} = [\hat{j}_1, \dots, \hat{j}_d]^T$ and $\tilde{\mathbf{j}} = [\tilde{j}_1, \dots, \tilde{j}_d]^T$ belonging in Θ , $\hat{\mathbf{j}}$ is ordered before $\tilde{\mathbf{j}}$ if for the smallest number k such that $\hat{j}_k \neq \tilde{j}_k$, we have $\hat{j}_k < \tilde{j}_k$. Then for any $\mathbf{j} \in \Theta$, $s(\Theta, \mathbf{j}) \in \{1, \dots, |\Theta|\}$ is defined to be the position of \mathbf{j} in the sorted array.*

A projection operator that takes tensor values over the observed indices are denoted by \mathbb{P}_Θ : for any d -th order tensor $\mathbf{X} \in \mathbb{R}^{n \times \dots \times n}$, $\mathbb{P}_\Theta(\mathbf{X}) := [p_1, \dots, p_{|\Theta|}]^T \in \mathbb{R}^{|\Theta|}$ where $p_{s(\Theta, \mathbf{j})} = \mathbf{X}(\mathbf{j}_1, \dots, \mathbf{j}_d)$ for all $\mathbf{j} \in \Theta$. Tensor recovery here is to find an approximation of the data tensor $\mathbf{X}_{\text{exact}}$ based on the entries over the observation indices, i.e., $\mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})$. Since it is assumed that $|\Theta| \ll d^n$, the cost for generating $\mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})$ is small compared with the cost for generating the whole $\mathbf{X}_{\text{exact}}$. When the CP rank of $\mathbf{X}_{\text{exact}}$ (denoted by R) is given, Acar et al. [17] formulate the tensor recovery problem as the following optimization problem

$$\begin{aligned}
&\min_{\mathbf{X}} \quad \frac{1}{2} \|\mathbb{P}_\Theta(\mathbf{X}) - \mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})\|_2^2 \\
&\text{s.t.} \quad \mathbf{X} = \sum_{r=1}^R \mathbf{v}_1^{(r)} \circ \mathbf{v}_2^{(r)} \circ \dots \circ \mathbf{v}_d^{(r)}, \tag{20}
\end{aligned}$$

It is clear that evaluating $\mathbb{P}_\Theta(\mathbf{X})$ requires $O(|\Theta|dR)$ flops.

To take the sparsity of gPC coefficients (5) into account (Cf. [6, 8] for sparse gPC approximations), a l_1 regularized version of (20) is formulated as follows [13]:

$$\min_{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d} J(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d) = \frac{1}{2} \|\mathbb{P}_\Theta([[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]]) - \mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})\|_2^2 + \beta \sum_{\|\mathbf{a}_i\|_1=0}^P |\langle [[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]], \mathbf{W}_i \rangle|, \tag{21}$$

where β is a regularization parameter, p is a given gPC order, and \mathbf{A}_k for $k = 1, \dots, d$ are CP factor matrices of \mathbf{X} (see (18)). To solve (21), the alternative minimization iterative method can be applied [13]. Letting $\mathbf{A}_k^{(q)}$ for $k = 1, \dots, d$ be the CP factor matrices at q -th iteration step ($q \geq 0$ is an integer), each CP factor matrix $\mathbf{A}_k^{(q+1)}$ at iteration step $q + 1$ is obtained through

$$\mathbf{A}_k^{(q+1)} = \arg \min_{\mathbf{A}_k} J(\mathbf{A}_1^{(q+1)}, \dots, \mathbf{A}_{k-1}^{(q+1)}, \mathbf{A}_k, \mathbf{A}_{k+1}^{(q)}, \dots, \mathbf{A}_d^{(q)}), \quad (22)$$

which leads to a generalized lasso problem and is discussed next.

3.3. A generalized lasso problem

Let $\text{vec}(\mathbf{A})$ denote the vector form of a given matrix \mathbf{A} (as implemented in the MATLAB function `reshape`). Following the procedures discussed in [13], (22) can be written as a generalized lasso problem

$$\text{vec}(\mathbf{A}_k^{(q+1)}) = \arg \min_s \frac{1}{2} \|\mathbf{B}s - \mathbf{b}\|_2^2 + \beta \|\mathbf{F}s\|_1, \quad (23)$$

where $\mathbf{B} := \mathcal{B}_{\Theta,k}([\mathbf{A}_1^{(q+1)}, \dots, \mathbf{A}_{k-1}^{(q+1)}, \mathbf{A}_k, \mathbf{A}_{k+1}^{(q)}, \dots, \mathbf{A}_d^{(q)}])$, $\mathbf{F} := \mathcal{F}_{\Upsilon,k}([\mathbf{A}_1^{(q+1)}, \dots, \mathbf{A}_{k-1}^{(q+1)}, \mathbf{A}_k, \mathbf{A}_{k+1}^{(q)}, \dots, \mathbf{A}_d^{(q)}])$ with $\mathcal{B}_{\Theta,k}(\cdot)$ and $\mathcal{F}_{\Upsilon,k}(\cdot)$ defined as follows (in Definition 2 and Definition 3 respectively), and $\mathbf{b} = \mathbb{P}_{\Theta}(\mathbf{X}_{\text{exact}}) \in \mathbb{R}^{|\Theta|}$.

Definition 2. For a given observation index set $\Theta \subset \mathbb{N}^d$ and $k \in \{1, \dots, d\}$, the operator $\mathcal{B}_{\Theta,k}$ defines a mapping: for a d -th order tensor $[[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]] \in \mathbb{R}^{n \times \dots \times n}$ with rank R , entries of $\mathbf{B} := \mathcal{B}_{\Theta,k}([\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]) \in \mathbb{R}^{|\Theta| \times nR}$ are zero except

$$\mathbf{B}(s(\Theta, \mathbf{j}), (r-1)n + j_k) = \prod_{k'=1}^{k-1} \mathbf{A}_{k'}(j_{k'}, r) \prod_{k'=k+1}^d \mathbf{A}_{k'}(j_{k'}, r), \text{ for all } \mathbf{j} = [j_1, \dots, j_d] \in \Theta \text{ and } r = 1, \dots, R, \quad (24)$$

where $s(\Theta, \cdot)$ is the sort operator defined in Definition 1.

Proposition 1. Let Θ be an observation index set, and $[[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]] \in \mathbb{R}^{n \times \dots \times n}$ be a d -th order tensor with rank R . Letting $\mathbf{A}_{i\delta} := [\mathbf{A}_i, \delta \mathbf{a}] \in \mathbb{R}^{n \times (R+1)}$ for $i = 1, 2, \dots, d$, for any $k \in \{1, \dots, d\}$, we have that

$$\mathcal{B}_{\Theta,k}([\mathbf{A}_{1\delta}, \mathbf{A}_{2\delta}, \dots, \mathbf{A}_{d\delta}]) = [\mathcal{B}_{\Theta,k}([\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]), \mathcal{B}_{\Theta,k}([\delta \mathbf{a}, \delta \mathbf{a}, \dots, \delta \mathbf{a}])]. \quad (25)$$

Proof. This proposition is straightforward since the matrix \mathbf{B} is constructed by (24). \square

Definition 3. Denoting the set of the multi-indices in (5) as $\Upsilon := \{\mathbf{i} | \mathbf{i} \in \mathbb{N}^d \text{ and } |\mathbf{i}| = 0, \dots, p\}$, where p is a given gPC order. The operator $\mathcal{F}_{\Upsilon,k}$ defines a mapping: for a d -th order tensor $[[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]] \in \mathbb{R}^{n \times \dots \times n}$ with rank R and $k \in \{1, \dots, d\}$, entries of $\mathbf{F} := \mathcal{F}_{\Upsilon,k}([\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d]) \in \mathbb{R}^{|\Upsilon| \times nR}$ are specified as

$$\mathbf{F}(s(\Upsilon, \mathbf{i}), (r-1)n + m_k) = \hat{\mathbf{w}}_k^{(i_k)}(m_k) \prod_{k'=1}^{k-1} \langle \mathbf{A}_{k'}(:, r), \hat{\mathbf{w}}_{k'}^{(i_{k'})} \rangle \prod_{k'=k+1}^d \langle \mathbf{A}_{k'}(:, r), \hat{\mathbf{w}}_{k'}^{(i_{k'})} \rangle, \quad (26)$$

for all $\mathbf{i} = [i_1, i_2, \dots, i_d]^T \in \Upsilon$, $m_k = 1, \dots, n$ and $r = 1, \dots, R$, where $s(\Upsilon, \cdot)$ is the sort operator defined in Definition 1 and each $\hat{\mathbf{w}}_k^{(i_k)}$ is defined in (12).

The operation numbers to construct \mathbf{B} and \mathbf{F} are $O(|\Theta|Rd)$ and $O(|\Upsilon|Rn^2d)$ respectively. The generalized lasso problem (23) can be solved by alternating direction method of multipliers (ADMM) [26] as follows. First, (23) is rewritten as

$$\begin{aligned} \min_{\mathbf{s}, \mathbf{t}} \quad & \frac{1}{2} \|\mathbf{B}\mathbf{s} - \mathbf{b}\|_2^2 + \beta \|\mathbf{t}\|_1 \\ \text{s.t.} \quad & \mathbf{F}\mathbf{s} = \mathbf{t}. \end{aligned} \quad (27)$$

The augmented Lagrangian function of (27) is given by

$$L(\mathbf{s}, \mathbf{t}, \mathbf{z}, \varrho) = \frac{1}{2} \|\mathbf{B}\mathbf{s} - \mathbf{b}\|_2^2 + \beta \|\mathbf{t}\|_1 + \langle \mathbf{z}, \mathbf{F}\mathbf{s} - \mathbf{t} \rangle + \frac{\varrho}{2} \|\mathbf{F}\mathbf{s} - \mathbf{t}\|_2^2,$$

where $\mathbf{z} \in \mathbb{R}^{|\Upsilon|}$ is the Lagrange multiplier, and $\varrho > 0$ is an augmented Lagrange multiplier. Following [27], the optimization problem (27) can be solved as

$$\begin{aligned} \mathbf{s}^{(i+1)} &= \arg \min_{\mathbf{s}} L(\mathbf{s}, \mathbf{t}^{(i)}, \mathbf{z}^{(i)}, \varrho^{(i)}) \\ \mathbf{t}^{(i+1)} &= \arg \min_{\mathbf{t}} L(\mathbf{s}^{(i+1)}, \mathbf{t}, \mathbf{z}^{(i)}, \varrho^{(i)}) \\ \mathbf{z}^{(i+1)} &= \mathbf{z}^{(i)} + \varrho^{(i)} (\mathbf{F}\mathbf{s}^{(i+1)} - \mathbf{t}^{(i+1)}). \end{aligned}$$

Following [26], details of the ADMM algorithm for (27) are summarized Algorithm 1, and the soft-thresholding operator $\mathbb{S}_{\frac{\beta}{\varrho}}$ on line 4 of Algorithm 1 is defined as

$$\mathbb{S}_{\frac{\beta}{\varrho}}(h) = \begin{cases} h - \frac{\beta}{\varrho} & \text{when } h \geq \frac{\beta}{\varrho}, \\ 0 & \text{when } |h| < \frac{\beta}{\varrho}, \\ h + \frac{\beta}{\varrho} & \text{when } h \leq -\frac{\beta}{\varrho}. \end{cases}$$

The cost of using Algorithm 1 to solve (27) is analyzed as follows:

- updating $\mathbf{s}^{(k+1)}$ line 3 of Algorithm 1 requires a matrix inversion and matrix-vector products, of which the total cost is $O(n^3R^3 + n^2R^2(|\Theta| + |\Upsilon|))$.
- updating the soft-thresholding operator and $\mathbf{z}^{(k+1)}$ requires $O(|\Upsilon|nR)$ operations.

Therefore, the total cost of Algorithm 1 is $C_{Alg1} := O(n^3R^3 + n^2R^2(|\Theta| + |\Upsilon|)) + O(|\Upsilon|Rn^2d + |\Theta|Rd)$.

The stopping criterion for the overall optimization problem (21) is specified through three parts in [13]: the relative changes of factor matrices, objective function values, and gPC coefficients. The relative change of factor matrices between iteration step q and $q + 1$ is defined as $\epsilon_{\text{factor}} := (\sum_{k=1}^d \|\mathbf{A}_k^{(q+1)} - \mathbf{A}_k^{(q)}\|_F^2)^{1/2} / (\sum_{k=1}^d \|\mathbf{A}_k^{(q)}\|_F^2)^{1/2}$ where $\|\cdot\|_F$ denotes the matrix Frobenius norm. The complexity of computing ϵ_{factor} is $O(dnR)$. Similarly, the relative changes of objective function values and gPC coefficients are defined as $\epsilon_J := |J^{(q+1)} - J^{(q)}|/|J^{(q)}|$ and $\epsilon_c := \|c^{(q+1)} - c^{(q)}\|_1 / \|c^{(q)}\|_1$ respectively, where $\|\cdot\|_1$ denotes the vector l_1 norm and $c^{(q)}$ collects the collocation coefficients (14) obtained with $[[\mathbf{A}_1^{(q)}, \dots, \mathbf{A}_d^{(q)}]]$. Since evaluating the objective function of (21) includes computing the projection \mathbb{P}_{Θ} and the tensor inner products, the cost of computing ϵ_J and ϵ_c are $O(|\Theta|dR + |\Upsilon|ndR)$. For a given tolerance δ , the optimization iteration for (21) terminates if $\epsilon_{\text{factor}} < \delta$, $\epsilon_J < \delta$ and $\epsilon_c < \delta$. The details for solving (21) is summarized Algorithm 2, which is proposed in [13]. The total cost of Algorithm 2 is $C_{Alg2} := dC_{Alg1} + O((n + |\Theta| + |\Upsilon|n)dR)$.

Algorithm 1 ADMM for generalized lasso [26]

Input: $\mathbf{B}, \mathbf{F}, \mathbf{b}, \beta, \nu \geq 1$ (for augment Lagrange multiplier)

- 1: Initialize $\varrho^{(0)}, \mathbf{s}^{(0)}, \mathbf{t}^{(0)} = \mathbf{F}\mathbf{s}^{(0)}$, and $\mathbf{z}^{(0)}, i = 0$
- 2: **while** not converged **do**
- 3: $\mathbf{s}^{(i+1)} = (\mathbf{B}^T \mathbf{B} + \varrho \mathbf{F}^T \mathbf{F})^\dagger (\mathbf{B}^T \mathbf{b} + \varrho \mathbf{F}^T \mathbf{t}^{(i)} - \mathbf{F}^T \mathbf{z}^{(i)})$
- 4: $\mathbf{t}^{(i+1)} = \mathbb{S}_{\frac{\beta}{\varrho}}(\mathbf{F}\mathbf{s}^{(i+1)} + \frac{1}{\varrho} \mathbf{z}^{(i)})$ (element-wise)
- 5: $\mathbf{z}^{(i+1)} = \mathbf{z}^{(i)} + \varrho(\mathbf{F}\mathbf{s}^{(i+1)} - \mathbf{t}^{(i+1)})$
- 6: $\varrho^{(i+1)} = \varrho^{(i)} \nu$
- 7: $i = i + 1$
- 8: **end while**

Output: $\mathbf{s}^* = \mathbf{s}^{(i)}$ and \mathbf{A} (the matrix form of \mathbf{s}^*).

Algorithm 2 Fixed-rank tensor recovery [13]

Input: CP rank R , initial rank R factor matrices $\mathbf{A}_k^{(0)}$ for $k = 1, \dots, d$, Θ , Υ and $\mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})$.

- 1: Let $q = 0$ and $\mathbf{b} = \mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})$.
- 2: Initialize $\epsilon_{\text{factor}} \geq \delta, \epsilon_J \geq \delta, \epsilon_c \geq \delta$.
- 3: **while** $\epsilon_{\text{factor}} \geq \delta, \epsilon_J \geq \delta$ or $\epsilon_c \geq \delta$ **do**
- 4: $q = q + 1$.
- 5: **for** $k = 1 : d$ **do**
- 6: $\mathbf{B} := \mathcal{B}_{\Theta,k}([\mathbf{A}_1^{(q+1)}, \dots, \mathbf{A}_{k-1}^{(q+1)}, \mathbf{A}_k, \mathbf{A}_{k+1}^{(q)}, \dots, \mathbf{A}_d^{(q)}])$.
- 7: $\mathbf{F} := \mathcal{F}_{\Upsilon,k}([\mathbf{A}_1^{(q+1)}, \dots, \mathbf{A}_{k-1}^{(q+1)}, \mathbf{A}_k, \mathbf{A}_{k+1}^{(q)}, \dots, \mathbf{A}_d^{(q)}])$.
- 8: Obtain factor matrix $\mathbf{A}_k^{(q+1)}$ by Algorithm 1.
- 9: **end for**
- 10: Compute $\epsilon_{\text{factor}}, \epsilon_J, \epsilon_c$.
- 11: **end while**

Output: CP factor matrices $\mathbf{A}_k = \mathbf{A}_k^{(q)}$ for $k = 1, 2, \dots, d$ and the recovered tensor $\mathbf{X} = [[\mathbf{A}_1, \dots, \mathbf{A}_d]]$.

4. Rank adaptive tensor recovery for stochastic collocation

Our goal is to perform uncertainty propagation from a high-dimensional random input vector ξ to the snapshot $\mathbf{y} = [u(x^{(1)}, \xi), \dots, u(x^{(N_h)}, \xi)]^T \in \mathbb{R}^{N_h}$ which is also high-dimensional. For this purpose, we develop a novel rank adaptive tensor recovery collocation (RATR-collocation) approach in this section. We first present our new general rank adaptive tensor recovery (RATR) algorithm for a general tensor, and then analyze its stability. After that, we present our main algorithm for this high-dimensional forward UQ problem.

4.1. Rank adaptive tensor recovery (RATR)

As discussed in section 3, the standard tensor recovery quadrature requires a given rank of the data tensor, which causes difficulties for problems where the tensor rank is not given a priori. Especially in our setting, tensor recovery quadratures are applied to compute the gPC coefficients for each kPCA mode (5), where the ranks of data tensors (10) are not given and the data tensor ranks associated with different kPCA modes can be different. To address this issue, we develop a new rank adaptive tensor recovery (RATR) approach.

Our idea is that, starting with setting the CP rank $R = 1$, we gradually increase the CP rank until the recovered tensor \mathbf{X} approximates the exact tensor $\mathbf{X}_{\text{exact}}$ well. To measure the quality of the recovered tensor, the following quantity of error is introduced

$$\varepsilon_{\Theta'}(\mathbf{X}) = \frac{\|\mathbb{P}_{\Theta'}(\mathbf{X}) - \mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}})\|_2}{\|\mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}})\|_2}, \quad (28)$$

where Θ is the observation index set, Θ' is a validation index set (see [28] for validation) randomly sampled from $\Theta_{\text{full}} := \{[j_1, j_2, \dots, j_d]^T \in \mathbb{N}^d \mid j_k = 1, \dots, n \text{ for } k = 1, \dots, d\}$, such that $\Theta' \cap \Theta = \emptyset$ and $|\Theta'| < |\Theta| \ll d^n$. Evaluating the relative error $\varepsilon_{\Theta'}(\mathbf{X})$ requires $O((|\Theta| + |\Theta'|)dR)$ flops, which is discussed in section 3.2.

Since the optimization problem (21) is non-convex, initial factor matrices in Algorithm 2 need to be properly chosen. We provide a detailed analysis of the initialization strategy in section 4.2. Here, supposing the tensor $\mathbf{X}^{(R)} = [[\mathbf{A}_1, \dots, \mathbf{A}_d]]$ is obtained, where $\{\mathbf{A}_1, \dots, \mathbf{A}_d\}$ is the solution of (21) with rank $R \geq 1$, we consider one higher rank, i.e., $R + 1$. While an analogous approach for tensor completion using tensor train decomposition can be found in [29], we here focus on CP decomposition and give the following scheme of rank-one update. The initial factor matrices for rank $R + 1$ is set to the rank-one updates of the factor matrices of $\mathbf{X}^{(R)}$, i.e.,

$$\mathbf{A}_k^{(0)} = [\mathbf{A}_k, \delta \mathbf{a}], \quad \text{for } k = 1, \dots, d, \quad (29)$$

where $\delta \mathbf{a} \in \mathbb{R}^n$ is a random perturbation vector and $\mathbf{A}_1, \dots, \mathbf{A}_d$ are the factor matrices of $\mathbf{X}^{(R)}$. With these new initial factor matrices, the recovered tensor $\mathbf{X}^{(R+1)}$ for rank $R + 1$ are obtained using Algorithm 2. To assess the progress obtained through this update of the CP rank, the difference between the recovery errors of $\mathbf{X}^{(R+1)}$ and $\mathbf{X}^{(R)}$ are assessed through $\Delta \varepsilon_{\Theta'} := \varepsilon_{\Theta'}(\mathbf{X}^{(R)}) - \varepsilon_{\Theta'}(\mathbf{X}^{(R+1)})$, where $\varepsilon_{\Theta'}(\cdot)$ is computed through (28). After that, we update the CP rank $R := R + 1$, and the above procedure is repeated until $\Delta \varepsilon_{\Theta'} < 0$, i.e., $\varepsilon_{\Theta'}(\mathbf{X}^{(R+1)}) > \varepsilon_{\Theta'}(\mathbf{X}^{(R)})$.

Details of our RATR method are presented in Algorithm 3. The initial rank-one matrices $\mathbf{A}_1^{(0)}, \dots, \mathbf{A}_d^{(0)}$ in the input are discussed in the next section. The other inputs are the observation index set Θ , the validation index set Θ' , and

entries of the data tensor (see (10)) on these index sets ($\mathbb{P}_\Theta(\mathbf{X}_{\text{exact}})$ and $\mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}})$). To start the **While** loop of this algorithm, $\Delta\epsilon_{\Theta'}$ is initially set to an arbitrary number that is larger than 0 on line 1. The output of this algorithm gives an estimation of the data tensor and its estimated rank. The cost of this algorithm is $C_{\text{Alg}}2 + O(|\Theta'|dR)$.

Algorithm 3 Rank adaptive tensor recovery (RATR)

Input: $\Theta, \Theta', \Upsilon, \mathbb{P}_\Theta(\mathbf{X}_{\text{exact}}), \mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}})$, and initial rank-one matrices $\mathbf{A}_k^{(0)}$ for $k = 1, \dots, d$.

- 1: Initialize the CP rank $R := 1$ and set $\Delta\epsilon_{\Theta'} > 0$.
- 2: Run Algorithm 2 to obtain $\mathbf{X}^{(1)}$.
- 3: Compute $\epsilon_{\Theta'}(\mathbf{X}^{(1)})$ by (28).
- 4: **while** $\Delta\epsilon_{\Theta'} \geq 0$ **do**
- 5: Initialize $\mathbf{A}_k^{(0)} := [\mathbf{A}_k, \delta\mathbf{a}]$ for $k = 1, \dots, d$, where each \mathbf{A}_k is a factor matrix of $\mathbf{X}^{(R)}$.
- 6: Update the CP rank $R := R + 1$.
- 7: Run Algorithm 2 to obtain $\mathbf{X}^{(R)}$.
- 8: Compute $\epsilon_{\Theta'}(\mathbf{X}^{(R)})$ by (28).
- 9: Compute the relative change in errors $\Delta\epsilon_{\Theta'} := \epsilon_{\Theta'}(\mathbf{X}^{(R-1)}) - \epsilon_{\Theta'}(\mathbf{X}^{(R)})$.
- 10: **end while**
- 11: Let $\mathbf{X} := \mathbf{X}^{(R-1)}$.
- 12: Let $R := R - 1$.

Output: the CP rank R and the recovered tensor \mathbf{X} .

4.2. Numerical stability analysis for RATR

While the tensor recovery problem (21) is a non-convex optimization problem, the initial guesses for the factor matrices need to be chosen properly. As discussed in section 3, (21) is solved using the alternative minimization iterative method, where the generalized lasso problem (23) needs to be solved at each iteration step. As studied in [30], (23) becomes ill-defined if \mathbf{B} is ill-conditioned. Therefore, a necessary condition for the initial factor matrices in (21) is that the resulting matrix \mathbf{B} (see Definition 2) needs to be well-conditioned. In this section, we first show that if the initial factor matrices are sampled through some given distributions, the condition number of \mathbf{B} is bounded with high probability for the case of rank $R = 1$. Next, we focus on the rank-one update procedure in our RATR approach (on line 5 of Algorithm 3), and show that the condition number of \mathbf{B} in (23) associated this update procedure is bounded under certain conditions. We begin our analysis with introducing the following definitions.

Definition 4 (Uniform observation index set). *An observation index set Θ is uniform if and only if $|\{j|\mathbf{j} = [j_1, \dots, j_d] \in \Theta \text{ and } j_k = i\}| = |\Theta|/n$ for each $i = 1, \dots, n$ and $k = 1, \dots, d$, where $|\cdot|$ denotes the size of a set.*

Definition 5 (d -th order ratio with m degrees of freedom). *Let $\psi_1, \dots, \psi_{d-1}$ form a random sample from a given distribution \mathbb{P} and for a given positive integer d , let $\Psi := \prod_{j=1}^{d-1} \psi_j$. For a given positive integer m , let Ψ_1, \dots, Ψ_m form a random sample from the sampling distribution of Ψ , and let $\Xi := \sum_{k=1}^m \Psi_k^2$. The d -th order ratio with m degrees of*

freedom of \mathbb{P} is

$$\mu = \frac{\sqrt{\text{Var}(\Xi)}}{\mathbb{E}(\Xi)}, \quad (30)$$

where $\mathbb{E}(\Xi)$ and $\text{Var}(\Xi)$ are the expectation and the variance of Ξ respectively.

Note that μ in (30) is typically called the coefficient of variation of Ξ [31, p. 845-846].

Theorem 1. Let $[[A_1, \dots, A_d]] \in \mathbb{R}^{n \times \dots \times n}$ be a d -th order rank-one tensor and Θ be an observation index set. Suppose that for $i = 1, \dots, d$, all entries of A_i form a random sample from distribution \mathbb{P} . Assume that the observation index set Θ is uniform. For a given constant $c_1 > 1$, if the d -th order ratio with $|\Theta|/n$ degrees of freedom of \mathbb{P} satisfies $\mu < 1$, $Q := (c_1 - 1)/(c_1\mu + \mu) > 1$, then the condition number of $\mathbf{B} = \mathcal{B}_{\Theta,k}([A_1, A_2, \dots, A_d])$ for $k = 1, 2, \dots, d$ (see Definition 2) satisfies $\text{cond}^2(\mathbf{B}) \leq c_1$ with probability at least $1 - 1/Q^2$.

Proof. Since the entries of \mathbf{B} are zero except $\mathbf{B}(s(\Theta, \mathbf{j}), j_k)$ for $\mathbf{j} = [j_1, \dots, j_d]^T \in \Theta$ (see Definition 2) where $s(\Theta, \cdot)$ is the sort operator, each row of \mathbf{B} can have at most one nonzero entry. Therefore, $\mathbf{B}^T \mathbf{B}$ is a diagonal matrix.

Since all entries of A_i form a random sample from distribution \mathbb{P} for $i = 1, \dots, d$, based on Definition 2, for $\mathbf{j} = [j_1, \dots, j_d]^T \in \Theta$, $k = 1, \dots, d$ and $j_k = 1, \dots, n$, we can express $\mathbf{B}(s(\Theta, \mathbf{j}), j_k) := \Psi := \prod_{j=1}^{d-1} \psi_j$, where $\psi_1, \dots, \psi_{d-1}$ form a random sample from \mathbb{P} . We next denote $\Psi_{s(\Theta, \mathbf{j}), j_k} := \mathbf{B}(s(\Theta, \mathbf{j}), j_k)$, where $\Psi_{s(\Theta, \mathbf{j}), j_k}$ form a random sample from the sampling distribution of Ψ for $\mathbf{j} = [j_1, \dots, j_d]^T \in \Theta$, $k = 1, \dots, d$ and $j_k = 1, \dots, n$. Therefore, with the assumption that Θ is uniform, we have $\mathbf{B}^T \mathbf{B} = \text{diag}(\Xi_1, \dots, \Xi_n)$ with $\Xi_i := \sum_{j=1}^{|\Theta|/n} (\Psi_j)^2$ for $i = 1, \dots, n$, where Ψ_j for $j = 1, \dots, |\Theta|/n$ form a random sample from the sampling distribution of Ψ .

According to the Chebyshev inequality,

$$\text{Prob}(|\Xi - \mathbb{E}(\Xi)| > Q \sqrt{\text{Var}(\Xi)}) \leq \frac{\text{Var}(\Xi)}{Q^2 \text{Var}(\Xi)},$$

which is equivalent to

$$\text{Prob}(\mathbb{E}(\Xi) - Q \sqrt{\text{Var}(\Xi)} \leq \Xi \leq \mathbb{E}(\Xi) + Q \sqrt{\text{Var}(\Xi)}) \geq 1 - \frac{1}{Q^2}. \quad (31)$$

Using $\text{diag}(\Xi_1, \dots, \Xi_n) = \mathbf{B}^T \mathbf{B}$ gives

$$\text{cond}^2(\mathbf{B}) = \frac{\max\{\Xi_1, \dots, \Xi_n\}}{\min\{\Xi_1, \dots, \Xi_n\}} \leq \frac{\mathbb{E}(\Xi) + Q \sqrt{\text{Var}(\Xi)}}{\mathbb{E}(\Xi) - Q \sqrt{\text{Var}(\Xi)}} = \frac{1 + Q\mu}{1 - Q\mu}. \quad (32)$$

Noting that $Q := (c_1 - 1)/(c_1\mu + \mu) > 1$ which gives $c_1 = (1 + Q\mu)/(1 - Q\mu)$, combining (31) and (32) establishes

$$\text{cond}^2(\mathbf{B}) \leq c_1$$

with probability at least $1 - 1/Q^2$. □

The conditions in Theorem 1 require that $\mu < 1$ (μ is defined in Definition 5) and imply $Q\mu < 1$, such that $\text{cond}(\mathbf{B})$ is bounded above with probability at least $1 - 1/Q^2$. To achieve a high probability for a bounded $\text{cond}(\mathbf{B})$, Q should be large and μ should then be small. So, the initial factor matrices (inputs of Algorithm 2) should be generated using realizations of a distribution \mathbb{P} of which μ is small. As an example, we show the estimated μ (the d -th order ratio with

Table 1: Examples of the d -th order ratio with $|m|$ degrees of freedom for several standard distributions with $d = 48$ and $m = 33$ and the corresponding values of $\text{cond}(\mathbf{B})$.

Distribution \mathbb{P}	Estimated μ	Average $\text{cond}(\mathbf{B})$
U(1, 2)	0.0394	1.9031
N(9, 0.1)	0.0242	1.1491
U(1, 3)	0.1051	3.7348
N(9, 0.5)	0.1390	1.1830
U(0, 1)	231.4262	2.27×10^3
N(0, 1)	163.5804	1.89×10^4

$|m|$ degrees of freedom) for several standard distributions in Table 1, where we set $d = 48$ and $m = 33$. Here, $U(a_1, a_2)$ refers to a uniform distribution on the interval $[a_1, a_2]$, and $N(a_1, a_2)$ refers to a normal distribution with mean a_1 and standard deviation a_2 . To compute the estimated μ in Table 1, $\mathbb{E}(\Xi)$ and $\text{Var}(\Xi)$ in (30) are computed using the sample mean and the sample variance of 10^5 samples of Ξ (note that the relationship between Ξ and \mathbb{P} is stated in Definition 5). In the procedure of generating each sample of Ξ , the corresponding \mathbf{B} is formulated and its condition number $\text{cond}(\mathbf{B})$ is stored (see Definition 5 and Theorem 1 for the relationship between Ξ and \mathbf{B}). The 10^5 samples of Ξ are associated with 10^5 samples of $\text{cond}(\mathbf{B})$, and Table 1 also shows the average of these samples of $\text{cond}(\mathbf{B})$ associated with each distribution. As shown in Table 1, the distributions listed above the dash line have $\mu < 1$, and they therefore can be used to generate initial factor matrices, while $U(0, 1)$ and $N(0, 1)$ should not be used.

Next, our analysis proceeds through induction. That is, supposing for a rank R tensor $\mathbf{X}^{(R)} = [[\mathbf{A}_1, \dots, \mathbf{A}_d]]$, its corresponding \mathbf{B} (see Definition 2) is well-conditioned, we show that the matrix \mathbf{B} associated with $\mathbf{X}^{(R+1)} = [[\mathbf{A}_1^{(0)}, \dots, \mathbf{A}_d^{(0)}]]$ is also well-conditioned, where $\mathbf{A}_k^{(0)} = [\mathbf{A}_k, \delta \mathbf{a}]$ for $k = 1, \dots, d$ are the rank-one updates of the factor matrices and $\delta \mathbf{a} \in \mathbb{R}^n$ is a perturbation vector. Before introducing our main theorem (Theorem 2), the following lemma is given.

Lemma 1. *Given two matrices $\mathbf{X}_1 \in \mathbb{R}^{n_1 \times n_2}$ and $\mathbf{X}_2 \in \mathbb{R}^{n_1 \times n_3}$ with full column ranks where $n_1 > n_2 \geq n_3$, let their singular value decompositions be $\mathbf{X}_1 = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T$ and $\mathbf{X}_2 = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}_2^T$, where $\mathbf{U}_1 \in \mathbb{R}^{n_1 \times n_1}$, $\mathbf{\Sigma}_1 \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{V}_1 \in \mathbb{R}^{n_2 \times n_2}$, $\mathbf{U}_2 \in \mathbb{R}^{n_1 \times n_1}$, $\mathbf{\Sigma}_2 \in \mathbb{R}^{n_1 \times n_3}$ and $\mathbf{V}_2 \in \mathbb{R}^{n_3 \times n_3}$, and let $\mathbf{u}_1^{(i)} \in \mathbb{R}^{n_1}$ for $i = 1, \dots, n_1$ and $\mathbf{u}_2^{(j)} \in \mathbb{R}^{n_1}$ for $j = 1, \dots, n_1$ denote the left singular vectors of \mathbf{X}_1 and \mathbf{X}_2 respectively. Assume the following two conditions hold: first*

$$\mathbf{\Sigma}_2 \mathbf{\Sigma}_2^T = \text{diag}(\underbrace{\lambda, \dots, \lambda}_{n_3}, \underbrace{0, \dots, 0}_{n_1 - n_3}), \quad (33)$$

where λ is a positive constant; second

$$\langle \mathbf{u}_2^{(j)}, \mathbf{u}_1^{(i)} \rangle = 0, \text{ for } j = n_3 + 1, \dots, n_1 \text{ and } i = 1, \dots, n_1. \quad (34)$$

Then $\mathbf{X}_1 \mathbf{X}_1^T$ and $\mathbf{X}_2 \mathbf{X}_2^T$ commute, i.e., $\mathbf{X}_1 \mathbf{X}_1^T \mathbf{X}_2 \mathbf{X}_2^T = \mathbf{X}_2 \mathbf{X}_2^T \mathbf{X}_1 \mathbf{X}_1^T$.

Proof. Let $\mathbf{P} = \mathbf{U}_2^T \mathbf{U}_1 \in \mathbb{R}^{n_1 \times n_1}$, $\mathbf{X} = \mathbf{P} \mathbf{\Sigma}_1 \mathbf{\Sigma}_1^T \mathbf{P}^T \in \mathbb{R}^{n_1 \times n_1}$ and $\mathbf{Z} = \mathbf{X} \mathbf{\Sigma}_2 \mathbf{\Sigma}_2^T \in \mathbb{R}^{n_1 \times n_1}$. Using (33), $\mathbf{Z}(i, j) = 0$ for $j = n_3 + 1, \dots, n_1$ and $i = 1, \dots, n_1$, while $\mathbf{Z}(i, j) = \lambda \mathbf{X}(i, j)$ for $i = 1, \dots, n_1$ and $j = 1, \dots, n_3$. Denoting the singular

values of $\mathbf{X}_1 \in \mathbb{R}^{n_1 \times n_2}$ by $\sigma_1, \dots, \sigma_{n_2}$ (note that $n_1 > n_2$) gives

$$\mathbf{X} = \sum_{i=1}^{n_2} \sigma_i^2 \mathbf{p}_i \mathbf{p}_i^T, \quad (35)$$

where \mathbf{p}_i for $i = 1, \dots, n_1$ are columns of \mathbf{P} . Since $\mathbf{p}_i = \mathbf{U}_2^T \mathbf{u}_1^{(i)}$ for $i = 1, \dots, n_1$, each element of \mathbf{p}_i is $\mathbf{p}_i(j) = \langle \mathbf{u}_2^{(j)}, \mathbf{u}_1^{(i)} \rangle$. Using (34) gives $\mathbf{p}_i(j) = 0$ for $j = n_3 + 1, \dots, n_1$ and $i = 1, \dots, n_1$. Therefore, (35) gives $\mathbf{X}(i, j) = 0$ for $j = n_3 + 1, \dots, n_1$ and $i = n_3 + 1, \dots, n_1$. In summary, each entry of \mathbf{Z} is

$$\mathbf{Z}(i, j) = \begin{cases} \lambda \mathbf{X}(i, j) & \text{for } i \leq n_3, j \leq n_3, \\ 0 & \text{otherwise.} \end{cases} \quad (36)$$

Similarly, each entry of $\mathbf{Z}' := \Sigma_2 \Sigma_2^T \mathbf{X} \in \mathbb{R}^{n_1 \times n_1}$ is

$$\mathbf{Z}'(i, j) = \begin{cases} \lambda \mathbf{X}(i, j) & \text{for } i \leq n_3, j \leq n_3, \\ 0 & \text{otherwise.} \end{cases} \quad (37)$$

Combing (36)–(37) gives $\mathbf{Z} = \mathbf{Z}'$, and thus $\mathbf{X} \Sigma_2 \Sigma_2^T = \Sigma_2 \Sigma_2^T \mathbf{X}$, which leads to

$$\mathbf{U}_2^T \mathbf{U}_1 \Sigma_1 \Sigma_1^T \mathbf{U}_1^T \mathbf{U}_2 \Sigma_2 \Sigma_2^T = \Sigma_2 \Sigma_2^T \mathbf{U}_2^T \mathbf{U}_1 \Sigma_1 \Sigma_1^T \mathbf{U}_1^T \mathbf{U}_2.$$

Left multiplying both sides of the above equation by \mathbf{U}_2 and right multiplying them by \mathbf{U}_2^T give

$$\mathbf{X}_1 \mathbf{X}_1^T \mathbf{X}_2 \mathbf{X}_2^T = \mathbf{X}_2 \mathbf{X}_2^T \mathbf{X}_1 \mathbf{X}_1^T.$$

□

Theorem 2. Let Θ be an observation index set, and $[[\mathbf{A}_1 \dots, \mathbf{A}_d]] \in \mathbb{R}^{n \times n \times \dots \times n}$ be a d -th order tensor with rank R . Suppose that $\mathbf{A}_{i_\delta} = [\mathbf{A}_i, \delta \mathbf{a}]$ with $\delta \mathbf{a} = [1, \dots, 1]^T \in \mathbb{R}^n$ is a rank-one update of \mathbf{A}_i for $i = 1, \dots, d$. Let $\mathbf{B} := \mathcal{B}_{\Theta, k}([[\mathbf{A}_1 \dots, \mathbf{A}_d]])$, $\delta \mathbf{B} := \mathcal{B}_{\Theta, k}([[\delta \mathbf{a}, \delta \mathbf{a}, \dots, \delta \mathbf{a}]])$, and $\mathbf{B}_\delta = \mathcal{B}_{\Theta, k}([[\mathbf{A}_{1\delta}, \dots, \mathbf{A}_{d\delta}]])$, and let $\mathbf{u}_{\delta \mathbf{B}}^{(j)}$ and $\mathbf{u}_{\mathbf{B}}^{(l)}$ are the j -th and the l -th left singular vectors of $\delta \mathbf{B}$ and \mathbf{B} respectively for $j, l = 1, \dots, |\Theta|$. If the following three conditions hold:

- 1) the observation index set Θ is uniform as defined in Definition 4,
- 2) there exist positive constants c_2 and c_3 which are independent of \mathbf{B} , such that $\text{cond}^2(\mathbf{B}) = s_{\max}^2(\mathbf{B})/s_{\min}^2(\mathbf{B}) \leq c_2$ and $|\Theta|/(ns_{\min}^2(\mathbf{B})) \leq c_3$, where s_{\max} and s_{\min} are the largest and the smallest singular values of \mathbf{B} respectively,
- 3) $\langle \mathbf{u}_{\delta \mathbf{B}}^{(j)}, \mathbf{u}_{\mathbf{B}}^{(l)} \rangle = 0$, for $j = n + 1, \dots, |\Theta|$ and $l = 1, \dots, |\Theta|$,

then the condition number of \mathbf{B}_δ satisfies that

$$\text{cond}^2(\mathbf{B}_\delta) \leq c_2 + c_3. \quad (38)$$

Proof. By Proposition 1, we have $\mathbf{B}_\delta = [\mathbf{B}, \delta \mathbf{B}]$. Note that the largest singular value of \mathbf{B}_δ satisfies that $s_{\max}^2(\mathbf{B}_\delta) = \lambda_{\max}(\mathbf{B} \mathbf{B}^T + \delta \mathbf{B} \delta \mathbf{B}^T)$. Using the min-max theorem,

$$\begin{aligned} s_{\max}^2(\mathbf{B}_\delta) &= \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^T (\mathbf{B} \mathbf{B}^T + \delta \mathbf{B} \delta \mathbf{B}^T) \mathbf{x} \\ &\leq \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^T (\mathbf{B} \mathbf{B}^T) \mathbf{x} + \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^T (\delta \mathbf{B} \delta \mathbf{B}^T) \mathbf{x} \\ &= s_{\max}^2(\mathbf{B}) + s_{\max}^2(\delta \mathbf{B}) \end{aligned}$$

Next we consider the smallest singular value of \mathbf{B}_δ under the above conditions. Let $\mathbf{B} = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T$ and $\delta \mathbf{B} = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}_2^T$ be the singular value decompositions of \mathbf{B} and $\delta \mathbf{B}$, respectively.

By condition 1) and noting that $\delta \mathbf{a} = [1, \dots, 1]^T$, we have $\mathbf{\Sigma}_2 \mathbf{\Sigma}_2^T = \text{diag}(\underbrace{|\Theta|/n, \dots, |\Theta|/n}_n, \underbrace{0, \dots, 0}_{|\Theta|-n})$. Since $\mathbf{\Sigma}_2 \mathbf{\Sigma}_2^T = \text{diag}(\underbrace{|\Theta|/n, \dots, |\Theta|/n}_n, \underbrace{0, \dots, 0}_{|\Theta|-n})$ and $\langle \mathbf{u}_{\delta \mathbf{B}}^{(j)}, \mathbf{u}_{\mathbf{B}}^{(l)} \rangle = 0$, $j = n+1, \dots, |\Theta|$, and $l = 1, \dots, |\Theta|$. By Lemma 1, $\mathbf{B} \mathbf{B}^T$ and $\delta \mathbf{B} \delta \mathbf{B}^T$ commute. Therefore, $\mathbf{B} \mathbf{B}^T$ and $\delta \mathbf{B} \delta \mathbf{B}^T$ can be simultaneously diagonalizable, i.e., there exists an orthonormal matrix \mathbf{Q} such that

$$\begin{aligned} \mathbf{B} \mathbf{B}^T &= \mathbf{Q} \mathbf{\Lambda}_1 \mathbf{Q}^T, \quad \mathbf{\Lambda}_1 = \text{diag}(s_{\max}^2(\mathbf{B}), \dots, s_{\min}^2(\mathbf{B}), \underbrace{0, \dots, 0}_{|\Theta|-n}), \\ \delta \mathbf{B} \delta \mathbf{B}^T &= \mathbf{Q} \mathbf{\Lambda}_2 \mathbf{Q}^T, \quad \mathbf{\Lambda}_2 = \text{diag}(s_{\max}^2(\delta \mathbf{B}), \dots, s_{\min}^2(\delta \mathbf{B}), \underbrace{0, \dots, 0}_{|\Theta|-n}). \end{aligned}$$

Then it follows that $\mathbf{B} \mathbf{B}^T + \delta \mathbf{B} \delta \mathbf{B}^T = \mathbf{Q}(\mathbf{\Lambda}_1 + \mathbf{\Lambda}_2) \mathbf{Q}^T$, and $s_{\min}^2(\mathbf{B}_\delta) = s_{\min}^2(\mathbf{B})$, which gives that

$$\begin{aligned} \text{cond}^2(\mathbf{B}_\delta) &= \frac{s_{\max}^2(\mathbf{B}_\delta)}{s_{\min}^2(\mathbf{B}_\delta)} \\ &\leq \frac{s_{\max}^2(\mathbf{B}) + s_{\max}^2(\delta \mathbf{B})}{s_{\min}^2(\mathbf{B})} = \text{cond}^2(\mathbf{B}) + |\Theta|/(n s_{\min}^2(\mathbf{B})) \leq c_2 + c_3. \end{aligned}$$

□

In summary, in our RATR algorithm (Algorithm 3), the fixed-rank tensor recovery algorithm (Algorithm 2) is invoked. The stability of Algorithm 2 is dependent on the observation index set Θ and the initial factor matrices. From our above analysis, if the observation index set Θ is uniform, and the initial rank-one factor matrices are sampled from the distributions given in Table 1 with $\mu < 1$, the first tensor recovery step in RATR (on line 2 of Algorithm 3) is stable with high probability. In the rank adaptive procedure, our analysis shows that the initial factor matrices specified on line 5 of Algorithm 3 can lead to stable tensor recovery on line 7 of Algorithm 3, if each \mathbf{B} (see Definition 1) associated with the data tensor obtained in the previous iteration step is well-conditioned. While the overall tensor recovery problem (21) is solved using the alternative minimization iterative method, our analysis is restricted to the first iteration step. To analyze the stability for the generalized lasso problem (23) for arbitrary iterations steps during the alternative minimization procedure remains an open problem. Nevertheless, our analysis here gives a systematic guidance to initialize the factor matrices for Algorithm 3 (also for Algorithm 2), and our numerical results in section 5 show that our RATR approach is stable and efficient.

4.3. RATR-collocation algorithm

Our goal is to efficiently conduct uncertainty propagation from the random input $\xi \in I^d$ to the discrete solution (which is high-dimensional) $\mathbf{y} = \chi(\xi) = [u(x^{(1)}, \xi), \dots, u(x^{(N_h)}, \xi)]^T \in \mathcal{M}$ of (1)–(2). The overall procedure of RATR-collocation approach is presented as the following three steps: generating data, processing data to construct RATR-collocation model, and conducting predictions using the RATR-collocation model.

For generating data, a tensor style quadrature rule [32] is first specified with n quadrature nodes in each dimension. The full index set is then defined as $\Theta_{\text{full}} := \{[j_1, j_2, \dots, j_d]^T \in \mathbb{N}^d \mid j_k = 1, \dots, n \text{ for } k = 1, \dots, d\}$ and quadrature

nodes are denoted as $\{\xi_{j_1 \dots j_d}, \text{ for } \mathbf{j} = [j_1, j_2, \dots, j_d]^T \in \Theta_{\text{full}}\}$. A observation index set Θ , and a validation index set Θ' are randomly selected from Θ_{full} , such that $\Theta' \cap \Theta = \emptyset$ and $|\Theta'| < |\Theta| \ll n^d$. After that, snapshots $\chi(\xi_{j_1 \dots j_d})$ for $\mathbf{j} = [j_1, j_2, \dots, j_d]^T \in \Theta \cup \Theta'$ are computed through solving deterministic versions of (1)–(2) with high-fidelity numerical schemes. At the end of this step, the snapshots are stored in a data matrix $\mathbf{Y} = [\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N_r)}]$, where $\mathbf{y}^{(s(\Theta \cup \Theta'), \mathbf{j})} := \chi(\xi_{j_1 \dots j_d})$ for $\mathbf{j} = [j_1, j_2, \dots, j_d]^T \in \Theta \cup \Theta'$ and $s(\cdot, \cdot)$ is the sort operator defined in Definition 1.

To process the data, kPCA (see section 2.1) is first applied to result in a reduced-dimensional representation of \mathbf{Y} —each $\mathbf{y}^{(l)} = \chi(\xi^{(l)}) \in \mathcal{M}$ is mapped to $\gamma(\xi^{(l)}) = [\gamma_1(\xi^{(l)}), \dots, \gamma_{N_r}(\xi^{(l)})]^T \in \mathcal{M}_r$ for $l = 1, \dots, N_r$. After that, for each kPCA mode $e = 1, \dots, N_r$, an estimated data tensor (10) is generated through our RATR approach presented in section 4.1. That is, through setting the observed data $\mathbb{P}_{\Theta}(\mathbf{X}_{\text{exact}}) := \mathbf{p}$ with $\mathbf{p} = [p_1, \dots, p_{|\Theta|}]^T$, where $p_{s(\Theta, \mathbf{j})} = \gamma_e(\xi^{(s(\Theta, \mathbf{j}))})$ for $\mathbf{j} \in \Theta$, and the validation data $\mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}}) := \mathbf{p}$ with $\mathbf{p} = [p_1, \dots, p_{|\Theta'|}]^T$, where $p_{s(\Theta', \mathbf{j})} = \gamma_e(\xi^{(s(\Theta', \mathbf{j}))})$ for $\mathbf{j} \in \Theta'$, Algorithm 3 gives an approximation of \mathbf{X}_e , which is denoted by $\tilde{\mathbf{X}}_e$. With this estimated data tensor, each gPC approximation (see (5)) $\gamma_e(\xi) \approx \gamma_e^{\text{gPC}} := \sum_{|\mathbf{i}|=0}^p c_{ei} \Phi_i(\xi)$ for $e = 1, \dots, N_r$ is obtained with coefficients computed through $c_{ei} := \langle \tilde{\mathbf{X}}_e, \mathbf{W}_i \rangle$, where \mathbf{W}_i is defined in (13). In the following, we call these gPC approximations $\{\gamma_e^{\text{gPC}}(\xi) := \sum_{|\mathbf{i}|=0}^p c_{ei} \Phi_i(\xi)\}_{e=1}^{N_r}$ the RATR-collocation model.

The above two steps for generating data and constructing the RATR-collocation model are summarized in Algorithm 4. For conducting a prediction of the snapshot for an arbitrary realization of ξ , we first use RATR-collocation model to compute the output $[\gamma_1^{\text{gPC}}(\xi), \dots, \gamma_{N_r}^{\text{gPC}}(\xi)]^T$ in the reduced-dimensional manifold \mathcal{M}_r . With the reduced output $[\gamma_1^{\text{gPC}}(\xi), \dots, \gamma_{N_r}^{\text{gPC}}(\xi)]^T$ and the data matrix \mathbf{Y} (generated in Algorithm 4), an estimation of the snapshot is obtained through the inverse mapping (see section 2.3 and [15]), which is denoted as $\mathbf{y}_{\text{RATR}} := \chi_{\text{RATR}}(\xi) \in \mathbb{R}^{N_h}$.

Algorithm 4 RATR-collocation in the reduced-dimensional manifold \mathcal{M}_r

Input: a full index set Θ_{full} , quadrature nodes $\{\xi_{j_1 \dots j_d}, \text{ for } \mathbf{j} = [j_1, j_2, \dots, j_d]^T \in \Theta_{\text{full}}\}$, an observation index set Θ , a validation index set Θ' , and a gPC order p .

- 1: Generate a data matrix $\mathbf{Y} = [\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N_r)}]$, where $\mathbf{y}^{(s(\Theta \cup \Theta'), \mathbf{j})} := \chi(\xi_{j_1 \dots j_d})$ for $\mathbf{j} \in \Theta \cup \Theta'$ are obtained through high-fidelity simulations for deterministic versions of (1)–(2) and $s(\cdot, \cdot)$ is defined in Definition 1.
- 2: Perform kPCA for \mathbf{Y} to obtain $\gamma(\xi^{(l)}) = [\gamma_1(\xi^{(l)}), \dots, \gamma_{N_r}(\xi^{(l)})]^T$ for $l = 1, \dots, N_r$.
- 3: **for** $e = 1 : N_r$ **do**
- 4: Define $\mathbb{P}_{\Theta}(\mathbf{X}_{\text{exact}}) := \mathbf{p}$ with $\mathbf{p} = [p_1, \dots, p_{|\Theta|}]^T$, where $p_{s(\Theta, \mathbf{j})} = \gamma_e(\xi^{(s(\Theta, \mathbf{j}))})$ for $\mathbf{j} \in \Theta$.
- 5: Define $\mathbb{P}_{\Theta'}(\mathbf{X}_{\text{exact}}) := \mathbf{p}$ with $\mathbf{p} = [p_1, \dots, p_{|\Theta'|}]^T$, where $p_{s(\Theta', \mathbf{j})} = \gamma_e(\xi^{(s(\Theta', \mathbf{j}))})$ for $\mathbf{j} \in \Theta'$.
- 6: Generate an estimated data tensor \mathbf{X} using Algorithm 3, and define $\tilde{\mathbf{X}}_e := \mathbf{X}$.
- 7: Generate the gPC approximation $\gamma_e(\xi) \approx \gamma_e^{\text{gPC}} := \sum_{|\mathbf{i}|=0}^p c_{ei} \Phi_i(\xi)$ with $c_{ei} := \langle \tilde{\mathbf{X}}_e, \mathbf{W}_i \rangle$ for $\mathbf{i} \in \Upsilon := \{\mathbf{i} | \mathbf{i} \in \mathbb{N}^d \text{ and } \|\mathbf{i}\|_1 = 0, \dots, p\}$, where \mathbf{W}_i is defined in (13).
- 8: **end for**

Output: gPC approximations $\gamma_1^{\text{gPC}}(\xi), \dots, \gamma_{N_r}^{\text{gPC}}(\xi)$ and the data matrix $\mathbf{Y} = [\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N_r)}]$.

5. Numerical study

In this section, we first consider diffusion problems in section 5.1 and section 5.2, and consider a Stokes problem in section 5.3. The governing equations of the diffusion problems are

$$-\nabla \cdot [a(x, \xi) \nabla u(x, \xi)] = 1 \quad \text{in} \quad D \times I^d, \quad (39)$$

$$u(x, \xi) = 0 \quad \text{on} \quad \partial D_D \times I^d, \quad (40)$$

$$\frac{\partial u(x, \xi)}{\partial n} = 0 \quad \text{on} \quad \partial D_N \times I^d, \quad (41)$$

where $\partial u / \partial n$ is the outward normal derivative of u on the boundaries, $\partial D_D \cap \partial D_N = \emptyset$ and $\partial D = \partial D_D \cup \partial D_N$. In the following numerical studies, the spatial domain is taken to be $D = (0, 1) \times (0, 1)$. The condition (40) is applied on the left ($x = 0$) and right ($x = 1$) boundaries, and (41) is applied on the top and bottom boundaries. Defining $H^1(D) := \{u : D \rightarrow \mathbb{R}, \int_D u^2 dD < \infty, \int_D (\partial u / \partial x_l)^2 dD < \infty, l = 1, 2\}$ and $H_0^1(D) := \{v \in H^1(D) | v = 0 \text{ on } \partial D_D\}$, the weak form of (39)–(41) is to find $u(x, \xi) \in H_0^1(D)$ such that $(a \nabla u, \nabla v) = (1, v)$ for all $v \in H_0^1(D)$. We discretize in space using a bilinear finite element approximation [33], with a uniform 65×65 grid ($N_h = 4225$).

The diffusion coefficient $a(x, \xi)$ in our numerical studies is assumed to be a random field with mean function $a_0(x)$, standard deviation σ and covariance function $Cov(x, y)$,

$$Cov(x, y) = \sigma^2 \exp\left(-\frac{|x_1 - y_1|}{l_c} - \frac{|x_2 - y_2|}{l_c}\right), \quad (42)$$

where $x = [x_1, x_2]^T, y = [y_1, y_2]^T \in \mathbb{R}^2$ and l_c is the correlation length. This random field can be approximated by a truncated Karhunen–Loève (KL) expansion [1]

$$a(x, \xi) \approx a_0(x) + \sum_{i=1}^d \sqrt{\lambda_i} a_i(x) \xi_i, \quad (43)$$

where $a_i(x)$ and λ_i are eigenfunctions and eigenvalues of (42), d is the number of KL modes retained, and $\{\xi_i\}_{i=1}^d$ are uncorrelated random variables. We set the random variables $\{\xi_i\}_{i=1}^d$ to be independent uniform distributions with range $I = [-1, 1]$, and set $a_0(x) = 1$ and $\sigma^2 = 0.25$. For test problem 1 (in section 5.1), we set $l_c = 0.8$ and $d = 48$, such that at least 95% of the total variance is captured, i.e., $(\sum_{i=1}^d \lambda_i) / (|D| \sigma^2) > 0.95$, where $|D|$ is the area of D . For test problem 2 (in section 5.2), we set $l_c = 1/16$ and again set $d = 48$.

For all test problems, we set the gPC order $p = 2$ (see section 2.2) and take $n = 3$ Gaussian quadrature points for each dimension, while Θ_{full} is constructed by the tensor product of these three points ($|\Theta_{\text{full}}| = 3^{48}$). As in the input of Algorithm 4, an observation index set Θ and a validation index Θ' are required. We test three cases of Θ uniformly sampled from Θ_{full} with sizes $|\Theta| = 100, 300$ and 600 respectively, and generate Θ' using 20 samples uniformly sampled from Θ_{full} , such that $\Theta \cap \Theta' = \emptyset$. Note that the number of high-fidelity simulations (the finite element methods here) in our RATR is $|\Theta| + |\Theta'|$, while that in standard tensor grid collocation [12] is $|\Theta_{\text{full}}| = 3^{48}$ and that in sparse grid collocation [3, 4] is still around 4705 (for a comparable grid level). So, the cost of RATR-collocation is much smaller than the costs of both tensor and sparse grid collocation methods for high-dimensional problems.

For the diffusion test problems. The regularization parameter β in (21) is set to 0.01, the tolerance in Algorithm 2 is set to $\delta = 10^{-5}$, and the initial rank-one matrices for Algorithm 3 are generated with samples of $U(1, 2)$ which

is an optimal initialization strategy as discussed in section 4.2. For kPCA as reviewed in section 2.1, we set the criterion for selecting principal components to $tol_{\text{PCA}} = 90\%$, and set the bandwidth to $\sigma_g = 5$ for the diffusion test problems. For a given realization of ξ , $\mathbf{y} := \chi(\xi)$ denotes the finite element solution, and $\mathbf{y}_{\text{RATR}} := \chi_{\text{RATR}}(\xi)$ refers to a RATR-collocation approximation solution (see section 4.3). A relative error is then defined as

$$\text{Relative error} = \frac{\|\mathbf{y} - \mathbf{y}_{\text{RATR}}\|_2}{\|\mathbf{y}\|_2}. \quad (44)$$

5.1. Test problem 1: diffusion problem with $l_c = 0.8$ and $d = 48$

For each case of the observation index set Θ (with $N_t := |\Theta| = 100, 300$ and 600 respectively), we first generate the corresponding data matrix \mathbf{Y} and apply kPCA for dimension reduction. For the given tolerance $tol_{\text{PCA}} = 90\%$, the number of kPCA modes retained is $N_r = 4$ for the three cases here (see section 2.1 for the definitions of N_r and tol_{PCA}). For each kPCA mode, our RATR algorithm gives an estimation $\tilde{\mathbf{X}}_e$ of the data tensor \mathbf{X}_e for $e = 1, \dots, N_r$ (see line 6 of Algorithm 4), where \mathbf{X}_e is defined in (10). Table 2 shows the estimated CP ranks of \mathbf{X}_e generated through Algorithm 3. It is clear that, these estimated ranks of each \mathbf{X}_e are similar for the three cases of Θ , and they are very small—the maximum estimated CP rank for this test problem is four.

Table 2: Estimated CP ranks of each data tensor \mathbf{X}_e for $e = 1, \dots, 4$, test problem 1.

rank \ e				
	1	2	3	4
$ \Theta $				
100	4	2	3	1
300	2	1	1	3
600	4	1	1	2

To assess the efficiency of our RATR procedure, we compare Algorithm 3 with the standard fixed-rank tensor recovery approach (Algorithm 2) to recover \mathbf{X}_1 with $|\Theta| = 600$ for this test problem. As discussed above, the initial rank-one factor matrices for RATR are generated through the distribution $U(1, 2)$. For Algorithm 2, for each given rank $R = 1, \dots, 4$, two distributions are tested for generating the initial matrices: $U(1, 2)$ and $N(0, 1)$. Note that, as discussed in section 4.2, $U(1, 2)$ is an optimal choice and $N(0, 1)$ is a non-optimal choice for the situation that the CP rank is one. In the following, the fixed-rank tensor recovery approach (Algorithm 2) with initial factor matrices generated through the optimal choice $U(1, 2)$ is denoted by FRTR-O, and that with initial factor matrices generated through the non-optimal choice $N(0, 1)$ is denoted by FRTR-N. Figure 2(a) shows the validation errors (28) of the recovered tensor generated by RATR, FRTR-O and FRTR-N respectively, where it is clear that for each rank $R = 1, \dots, 4$, our RATR has the smallest validation error. As discussed in section 3.2, the overall tensor recovery problem (21) is solved through the alternative minimization iterative method (see (22)). Looking more closely, the validation errors at each iteration step of the alternative minimization iterative method for $R = 1, 2, 4$ are shown in Figure 2(b), Figure 2(c) and Figure 2(d) respectively (since the results of $R = 3$ and $R = 4$ similar, we only show the results of

$R = 4$). For $R = 1$ (Figure 2(b)), there is no rank adaptive procedure preformed in RATR, and the validation errors of RATR and FRTR-O are the same, while it is clear that they are much smaller than the errors of FRTR-N. Moreover, the validation error of FRTR-N can even become larger as the iteration step increases for $R = 1$, which is consistent with Theorem 1. For $R = 2, 4$ (Figure 2(c) and Figure 2(d)), it can be seen that RATR has the smallest validation errors at each iteration step, which shows that our rank-one updating procedure (on line 5 of Algorithm 3) gives efficient initial factor matrices for the generalized lasso problem (23).

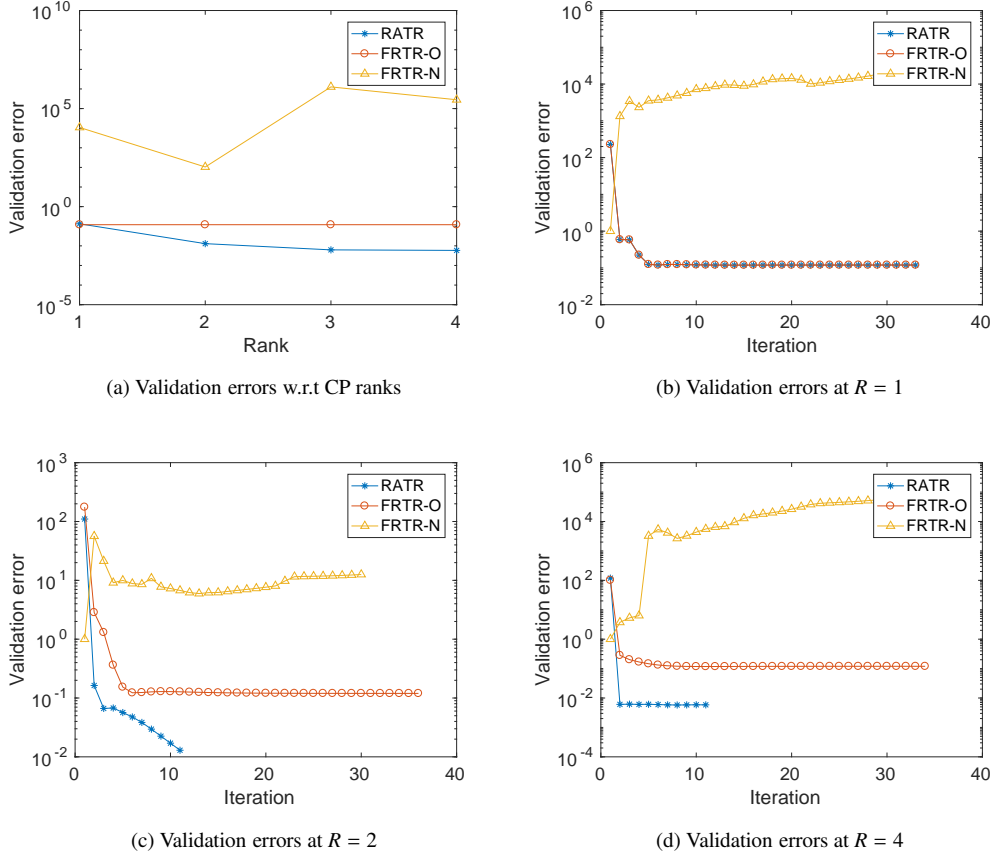


Figure 2: Validation errors of rank adaptive tensor recovery (RATR), fixed-rank tensor recovery with initial factor matrices generated through $U(1, 2)$ (FRTR-O), and fixed-rank tensor recovery with initial factor matrices generated through $N(0, 1)$ (FRTR-N), test problem 1.

While the sparsity of the gPC coefficients is taken into account in the tensor recovery problem (21), we show the absolute value of each the gPC coefficient c_{ei} (see (6)) for each kPCA mode $e = 1, \dots, 4$ and each gPC multi-index $i \in \Upsilon$ (see section 2.2) in Figure 3. In Figure 3, the gPC multi-index set is labeled as $\Upsilon = \{i^{(1)}, \dots, i^{|\Upsilon|}\}$, where the indices are sorted by the sort operator $s(\Upsilon, \cdot)$ (see section 3.2, Definition 1). From Figure 3, it is clear that the gPC coefficients are sparse—absolute values of most coefficients are smaller than 10^{-4} , which is consist with the results in [6].

Figure 4 shows the finite element solution \mathbf{y} and the RATR-collocation approximation \mathbf{y}_{RATR} responding to a given realization of ξ , where it can be seen that they are visually indistinguishable. Finally, we generate 500 samples of ξ , and compute the relative error (44) for the three cases ($|\Theta| = 100, 300$ and 600 respectively). Figure 5 shows Tukey

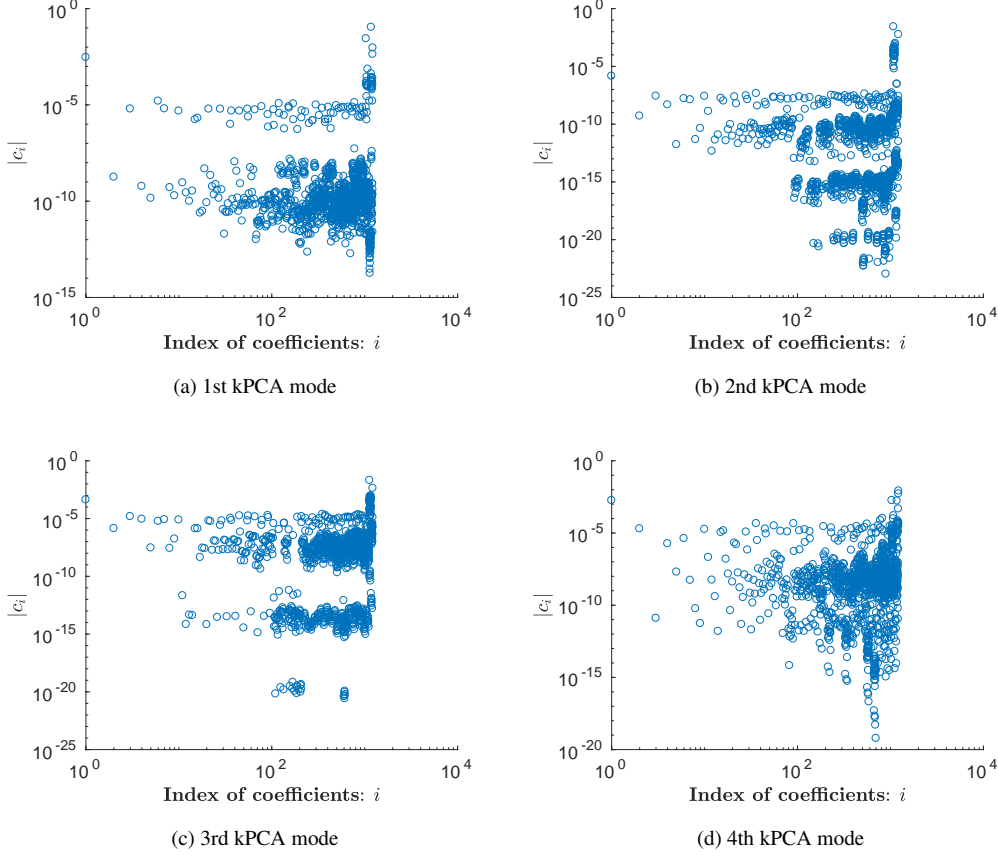
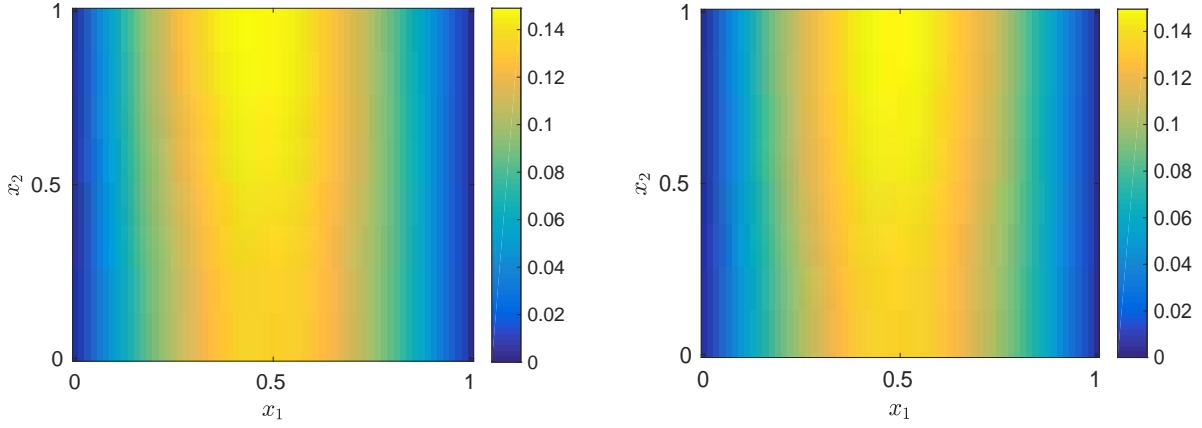


Figure 3: Sparsity of gPC coefficients for each kPCA mode, test problem 1.

box plots of these errors. Here, the central line in each box is the median, the lower and the upper edges are the first and the third quartiles respectively, and the red crosses are the outliers where the relative errors are large. From Figure 5, it is clear that as the size of the observation index set ($|\Theta|$) increases, values of the median, the first and the third quartiles of the errors decrease.

5.2. Test problem 2: diffusion problem with $l_c = 1/16$ and $d = 48$

For this test problem, the correlation length is very small, and the diffusion problem becomes highly non-smooth. Following the discussion procedure in test problem 1, we first generate the corresponding data matrix \mathbf{Y} for the three cases of Θ ($|\Theta| = 100, 300$ and 600) and apply kPCA on it. For $tol_{PCA} = 90\%$, the number of kPCA modes retained is $N_r = 7$ for this test problem. Tabel 3 shows the estimated CP ranks of \mathbf{X}_e generated through Algorithm 3 for each $e = 1, \dots, 7$, where it is clear that the estimated ranks are small (the maximum of the estimated ranks is seven). Figure 6 shows validation errors of our RATR, FRTR-O and FRTR-N (Algorithm 2 with initial factor matrices generated through $U(1, 2)$ and $N(0, 1)$ respectively) for recovering \mathbf{X}_1 (see (10)) with $|\Theta| = 600$. From Figure 6(a), it can be seen that our RATR has the smallest validation error for each rank $R = 2, \dots, 7$. It is also clear that, as the ranks increase, the error of RATR decreases, while the errors of FRTR-O and FRTR-N do not decrease. The other pictures in Figure 6 show the validation errors at each iteration step of the alternative minimization iterative procedure for $R = 1, 2, 3, 4, 7$



(a) Finite element solution

(b) RATR-collocation approximation

Figure 4: The finite element solution and the RATR-collocation approximation (with $|\Theta| = 600$) responding to a given realization of ξ , test problem 1.

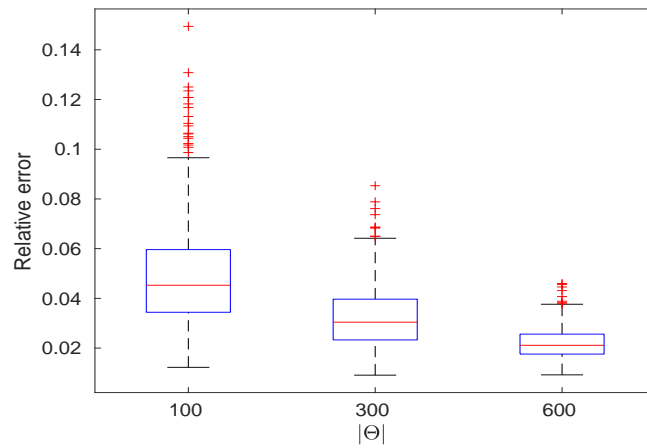


Figure 5: Relative errors of RATR-collocation approximation for 500 test samples, test problem 1.

(since the errors for $R = 5, 6$ are similar to the errors for $R = 7$, they are not shown here). For the case $R = 1$ (Figure 6(b)), while RATR is the same as FRTR-O, the error of RATR is larger than the error of FRTR-N, but the errors of RATR and FRTR-N are both very large (larger than one), which implies that this rank ($R=1$) is too small to accurately recover the data tensor. As the rank increases, for $R = 2, 3, 4, 7$, the validation errors of RATR are clearly smaller than the errors of FRTR-O and FRTR-N, which is consistent with the results in test problem 1. Figure 7 shows the absolute values of the gPC coefficient c_{ei} (see (6)) for $i \in \Upsilon$ and $e = 1, \dots, 4$ (the first four kPCA modes). It is clear that absolute values of most gPC coefficients are very small. Therefore, the gPC expansions for these four kPCA modes are sparse. For the other kPCA modes ($e = 5, 6, 7$), since the situation is similar to that of the first four kPCA modes, their corresponding gPC coefficients are not shown here, while these gPC expansions are also sparse. Finally, Figure 8 shows Tukey box plots of the relative errors (44) for 500 test samples for this test problem, where the central line in each box is the median, the lower and the upper edges are the first and the third quartiles respectively, and the red crosses are the outliers. From Figure 8, it is clear that, as the size of the observation index set ($|\Theta|$) increases, values of the median, the first and the third quartiles of the errors decrease, which are all consistent with the results in test problem 1.

Table 3: Estimated CP ranks of each data tensor \mathbf{X}_e for $e = 1, \dots, 7$, test problem 2.

rank \ e	1	2	3	4	5	6	7
$ \Theta $							
100	7	1	3	2	1	2	1
300	6	3	1	1	7	1	1
600	7	1	4	3	7	1	3

5.3. Test problem 3: the Stokes equations

The governing equations for this test problem are

$$\nabla \cdot [a(x, \xi) \nabla u(x, \xi)] + \nabla p(x, \xi) = 0 \quad \text{in } D \times I^d, \quad (45)$$

$$\nabla \cdot u(x, \xi) = 0 \quad \text{in } D \times I^d, \quad (46)$$

$$u(x, \xi) = g(x) \quad \text{on } \partial D \times I^d, \quad (47)$$

where $D \subset \mathbb{R}^2$, and $u(x, \xi) = [u_1(x, \xi), u_2(x, \xi)]^T$ and $p(x, \xi)$ are the flow velocity and the scalar pressure respectively. We consider the problem with uncertain viscosity $a(x, \xi)$, which is assumed to be a random field with mean function $a_0(x) = 1$, variance $\sigma^2 = 0.25$, and covariance function (42). The correlation length is set to $l_c = 0.8$, and we take $d = 48$ to capture 95% of the total variance as in test problem 1. We here consider the driven cavity flow problem posed on $D = (0, 1) \times (0, 1)$. For boundary conditions, the velocity profile $u = [1, 0]^T$ is imposed on the top boundary ($x_2 = 1$ where $x = [x_1, x_2]^T$), and $u = [0, 0]^T$ is imposed on all other boundaries. We discretize in space using the inf-sup stable $\mathcal{Q}_2 - \mathcal{P}_{-1}$ mixed finite element method (biquadratic velocity–linear discontinuous pressure [33]) as

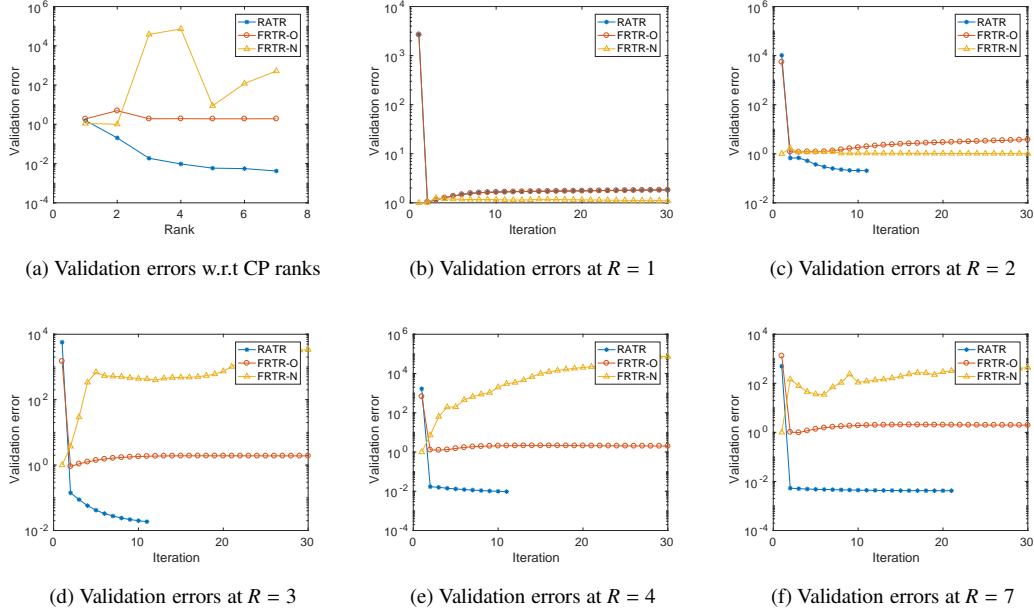


Figure 6: Validation errors of rank adaptive tensor recovery (RATR), fixed-rank tensor recovery with initial factor matrices generated through $U(1, 2)$ (FRTR-O), and fixed-rank tensor recovery with initial factor matrices generated through $N(0, 1)$ (FRTR-N), test problem 2.

implemented in IFISS [34] with a uniform 33×33 grid, which yields the velocity degrees of freedom $N_{h,u} = 2178$ and the pressure degrees of freedom $N_{h,p} = 768$. The output \mathbf{y} here is defined to be a vector collecting both discrete velocity and pressure solutions, and the overall dimension of \mathbf{y} is then $N_h = N_{h,u} + N_{h,p} = 2946$. For this test problem, the regularization parameter β in (21) is set to 0.1, and the tolerance in Algorithm 2 is set to $\delta = 10^{-5}$. The bandwidth σ_g of kPCA for dimension reduction is set to 10, and we again set tol_{PCA} to 90%.

We first generate the corresponding data matrix \mathbf{Y} for the three cases of Θ ($|\Theta| = 100, 300$ and 600) and apply kPCA on it. For $tol_{PCA} = 90\%$, our results show that the number of kPCA modes retained is $N_r = 9$ for the case $|\Theta| = 100$, while $N_r = 10$ for the cases $|\Theta| = 300$ and $|\Theta| = 600$, which implies that the sample size 100 may not be large enough for an accurate dimension reduction. Tabel 4 shows the estimated CP ranks of \mathbf{X}_e generated through Algorithm 3 for each kPCA mode $e = 1, \dots, N_r$. Again, it is clear that, the estimated ranks are small—the maximum estimated rank is only seven. This shows that the rank-one update produced in Algorithm 3 is performed seven times at most, and it is therefore not costly.

Figure 9 shows the validation errors of our RATR (Algorithm 3 with initial factor matrices generated through $U(1, 2)$), FRTR-O and FRTR-N (Algorithm 2 with initial factor matrices generated through $U(1, 2)$ and $N(0, 1)$ respectively) for recovering \mathbf{X}_1 (see (10)) with $|\Theta| = 600$. From Figure 9(a), it can be seen that our RATR has the smallest validation error for each rank. It is also clear that, as the rank increases from one to three, the errors of RATR reduces significantly, while the iterations from rank three to seven are caused by our stopping criterion on line 9 of Algorithm 3. The other pictures in Figure 9 show the validation errors at each iteration step of the alternative minimization iterative procedure for $R = 1, 2, 3, 6, 7$ (while the errors for $R = 4, 5$ are similar to the errors for $R = 6$, they are not shown here). Similarly to test problem 1, for the case $R = 1$, RATR is the same as FRTR-O, and their errors are smaller

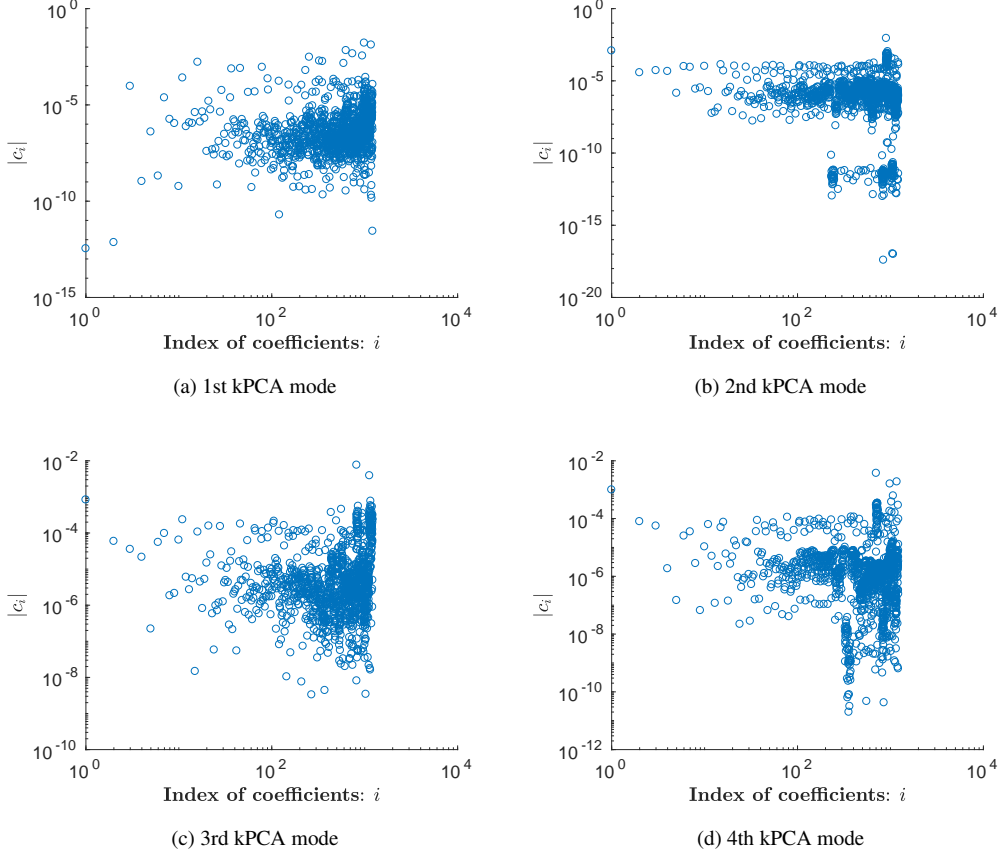


Figure 7: Sparsity of gPC coefficients for each kPCA mode, test problem 2.

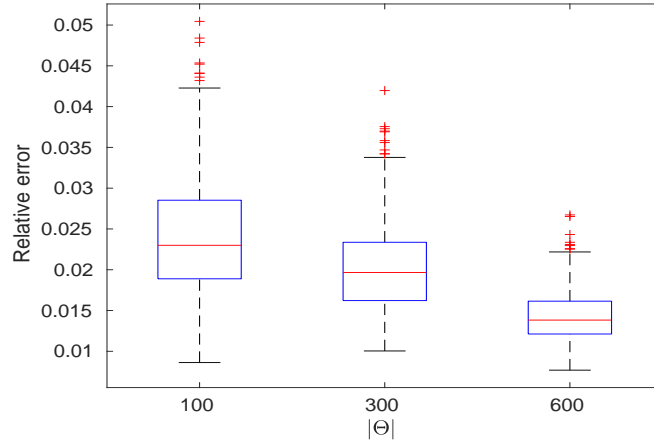


Figure 8: Relative errors of RATR-collocation approximation for 500 test samples, test problem 2.

than the error of FRTR-N. For $R = 2, 3, 6, 7$, the validation error of RATR is again clearly smaller than the errors of FRTR-O and FRTR-N, which shows that our rank-one update procedure is efficient for this Stokes problem.

Figure 10 shows the absolute values of the gPC coefficients of the first four kPCA modes for this test problem. It

is clear that absolute values of most gPC coefficients are small, and the gPC expansions for these four kPCA modes are therefore sparse. For the other kPCA modes ($e = 5, 6, 7$), while the situation is similar (the corresponding gPC expansions are also clearly sparse), they are not shown here. Figure 11 shows the flow streamlines and the pressure fields generated by the mixed finite element method and RATR-collocation (see section 4.3) responding to a given realization of ξ . It can be seen that there is no visual difference between the results obtained through finite elements and RATR-collocation. Finally, we generate 500 samples of ξ and compute the relative errors (44). Figure 8 shows Tukey box plots of these errors, where the central line in each box is the median, the lower and the upper edges are the first and the third quartiles respectively. It is clear that, as the size of the observation index set ($|\Theta|$) increases, values of the median, the first and the third quartiles of the errors decrease, which are all consistent with the results of the diffusion test problems.

Table 4: Estimated CP ranks of each data tensor \mathbf{X}_e for $e = 1, \dots, 10$, test problem 3.

rank e										
	1	2	3	4	5	6	7	8	9	10
$ \Theta $										
100	3	5	7	1	2	2	4	1	3	–
300	7	3	2	1	1	5	5	2	2	2
600	7	7	7	4	2	2	2	1	4	7

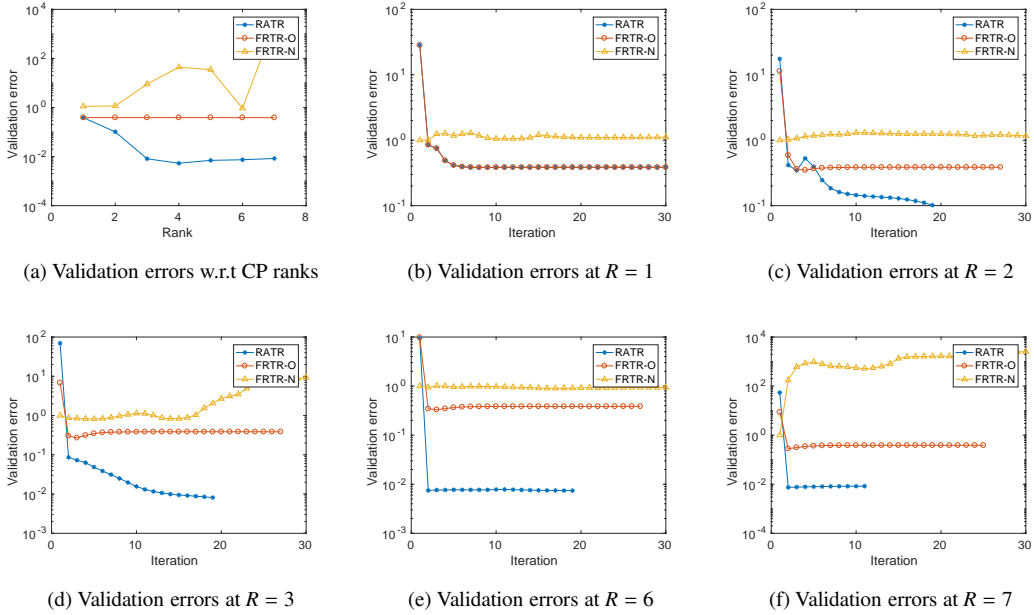


Figure 9: Validation errors of rank adaptive tensor recovery (RATR), fixed-rank tensor recovery with initial factor matrices generated through $U(1, 2)$ (FRTR-O), and fixed-rank tensor recovery with initial factor matrices generated through $N(0, 1)$ (FRTR-N), test problem 3.

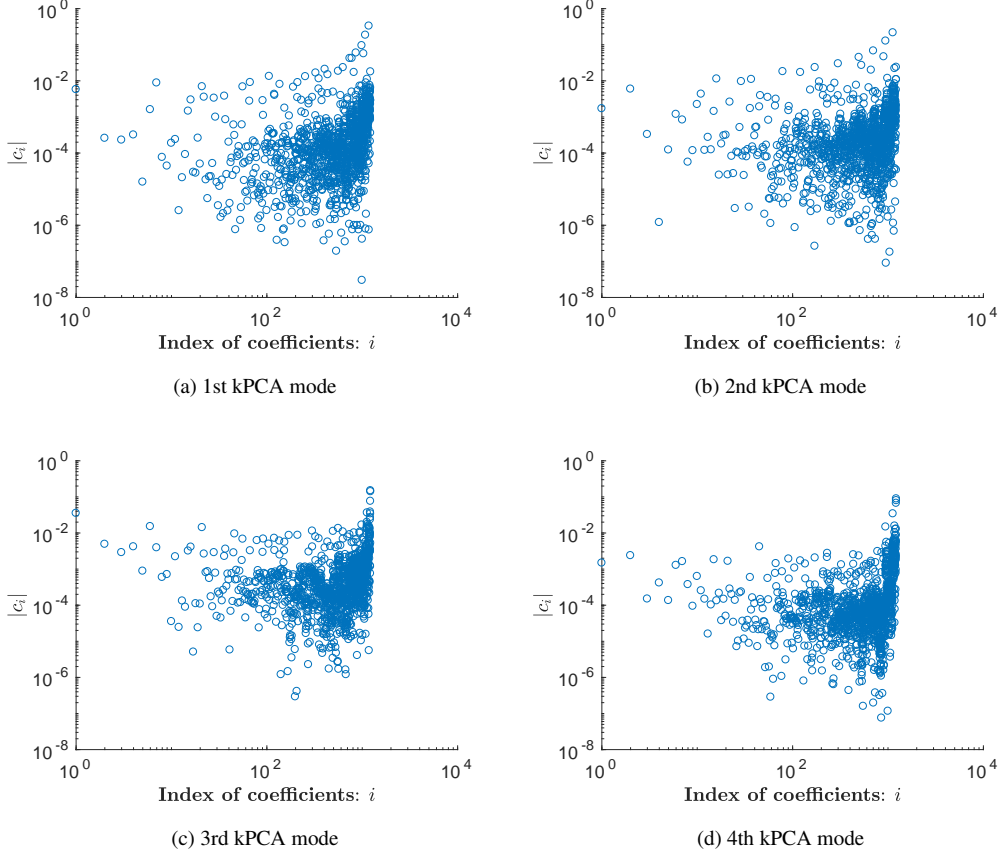
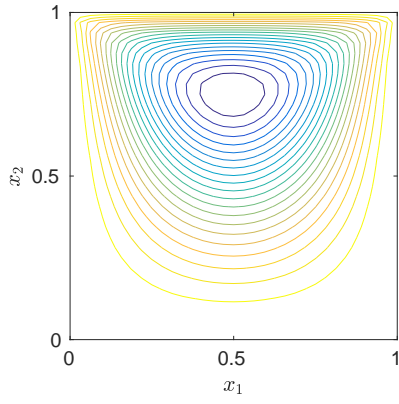


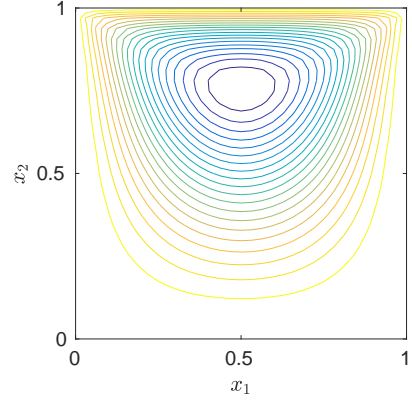
Figure 10: Sparsity of gPC coefficients for each kPCA mode, test problem 3.

6. Conclusions

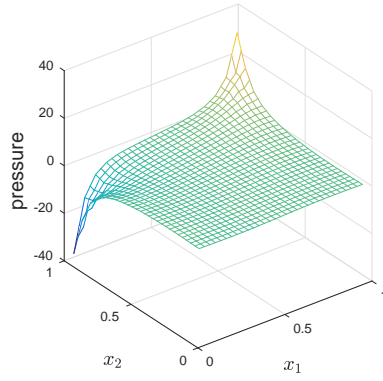
Exploiting potential low-dimensional structures is a fundamental concept of efficient surrogate and reduced order modelling for high-dimensional UQ problems. With a focus on the tensor recovery based stochastic collocation, our main conclusion is that our rank adaptive tensor recovery collocation (RATR-collocation) approach can efficiently exploit low-dimensional structures in this challenging problem in two aspects: first, we reformulate stochastic collocation based on manifold learning, where nonlinear low-dimensional structures in the snapshots are captured through kPCA; second, our novel RATR algorithm automatically explores the low-rank structures in the data tensors for computing the collocation coefficients without requiring a given tensor rank. Moreover, another main contribution of this work is the analysis of RATR, where the stability of our initialization strategies and the rank-one update procedure for the non-convex optimization problems involved is proven theoretically, such that a systematic guidance to initialize the factor matrices is provided to result in efficient and stable recovery results. As the performance of RATR algorithm depends on the CP rank of the data tensor (although it does not need to be explicitly given), our RATR-collocation is efficient when the CP rank is small, while it may not be efficient for high-rank problems. A possible solution for efficiently recovering high-rank tensors is to conduct adaptivity with respect to physical properties of the underlying PDE models, e.g., domain decomposition methods. Designing and analyzing such strategies will be the focus of our



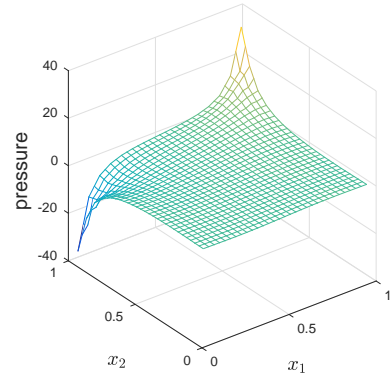
(a) Streamline: finite element solution



(b) Streamline: RATR-collocation approximation



(c) Pressure: finite element solution



(d) Pressure: RATR-collocation approximation

Figure 11: The finite element solution and the RATR-collocation approximation (with $|\Theta| = 600$) responding to a given realization of ξ , test problem 3.

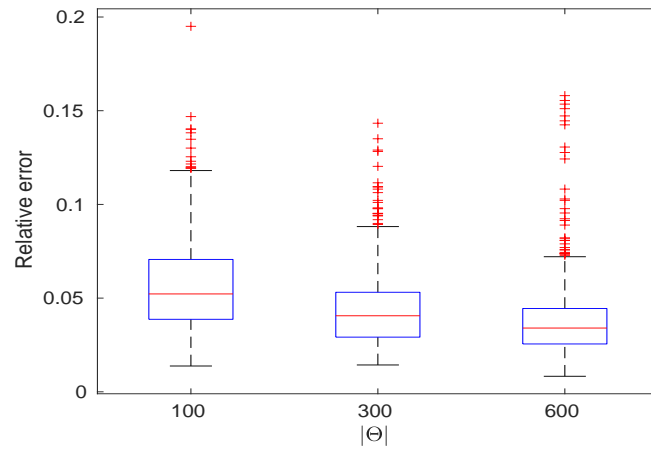


Figure 12: Relative errors of RATR-collocation approximation for 500 test samples, test problem 3.

future work.

Acknowledgments: This work is supported by the National Natural Science Foundation of China (No. 11601329) and the science challenge project (No. TZ2018001).

References

References

- [1] R. G. Ghanem, P. D. Spanos, *Stochastic Finite Elements: A Spectral Approach*, Courier Corporation, 2003.
- [2] D. Xiu, G. E. Karniadakis, The wiener–askey polynomial chaos for stochastic differential equations, *SIAM journal on scientific computing* 24 (2) (2002) 619–644.
- [3] D. Xiu, J. S. Hesthaven, High-order collocation methods for differential equations with random inputs, *SIAM Journal on Scientific Computing* 27 (3) (2005) 1118–1139.
- [4] X. Ma, N. Zabarar, An adaptive hierarchical sparse grid collocation algorithm for the solution of stochastic differential equations, *Journal of Computational Physics* 228 (8) (2009) 3084–3113.
- [5] C. Powell, H. Elman, Block-diagonal preconditioning for spectral stochastic finite-element systems, *IMA Journal of Numerical Analysis* 29 (2009) 350–375.
- [6] A. Doostan, H. Owhadi, A non-adapted sparse approximation of pdes with stochastic inputs, *Journal of Computational Physics* 230 (8) (2011) 3015–3034.
- [7] J. Peng, J. Hampton, A. Doostan, A weighted l_1 -minimization approach for sparse polynomial chaos expansions, *Journal of Computational Physics* 267 (2014) 92–111.
- [8] L. Yan, L. Guo, D. Xiu, Stochastic collocation algorithms using l_1 -minimization, *International Journal for Uncertainty Quantification* 2 (3).
- [9] J. D. Jakeman, A. Narayan, T. Zhou, A generalized sampling and preconditioning scheme for sparse approximation of polynomial chaos expansions, *SIAM Journal on Scientific Computing* 39 (3) (2017) A1114–A1144.
- [10] H. Lei, X. Yang, B. Zheng, G. Lin, N. A. Baker, Constructing surrogate models of complex systems with enhanced sparsity: quantifying the influence of conformational uncertainty in biomolecular solvation, *Multiscale Modeling & Simulation* 13 (4) (2015) 1327–1353.
- [11] L. Guo, A. Narayan, T. Zhou, A gradient enhanced l_1 -minimization for sparse approximation of polynomial chaos expansions, *Journal of Computational Physics* 367 (2018) 49–64.
- [12] I. Babuška, F. Nobile, R. Tempone, A stochastic collocation method for elliptic partial differential equations with random input data, *SIAM Journal on Numerical Analysis* 45 (3) (2007) 1005–1034.

- [13] Z. Zhang, T.-W. Weng, L. Daniel, Big-data tensor recovery for high-dimensional uncertainty quantification of process variations, *IEEE Transactions on Components, Packaging and Manufacturing Technology* 7 (5) (2017) 687–697.
- [14] S. Conti, A. OHagan, Bayesian emulation of complex multi-output and dynamic computer models, *Journal of statistical planning and inference* 140 (3) (2010) 640–651.
- [15] W. Xing, V. Triantafyllidis, A. Shah, P. Nair, N. Zabaras, Manifold learning for the emulation of spatial fields from computational models, *Journal of Computational Physics* 326 (2016) 666–690.
- [16] X. Ma, N. Zabaras, Kernel principal component analysis for stochastic input model generation, *Journal of Computational Physics* 230 (19) (2011) 7311–7331.
- [17] E. Acar, D. M. Dunlavy, T. G. Kolda, M. Mrup, Scalable tensor factorizations for incomplete data, *Chemometrics & Intelligent Laboratory Systems* 106 (1) (2010) 41–56.
- [18] S. Gandy, B. Recht, I. Yamada, Tensor completion and low-n-rank tensor recovery via convex optimization, *Inverse Problems* 27 (2) (2011) 025010.
- [19] J. Liu, P. Musialski, P. Wonka, J. Ye, Tensor completion for estimating missing values in visual data., in: *IEEE International Conference on Computer Vision*, 2013, pp. 2114–2121.
- [20] R. Vidal, Y. Ma, S. S. Sastry, Generalized principal component analysis, Vol. 5, Springer, 2016.
- [21] J.-Y. Kwok, I.-H. Tsang, The pre-image problem in kernel methods, *IEEE transactions on neural networks* 15 (6) (2004) 1517–1525.
- [22] B. Schölkopf, A. Smola, K.-R. Müller, Nonlinear component analysis as a kernel eigenvalue problem, *Neural computation* 10 (5) (1998) 1299–1319.
- [23] W. Gautschi, On generating orthogonal polynomials, *SIAM Journal on Scientific and Statistical Computing* 3 (3) (1982) 289–317.
- [24] L. D. Lathauwer, B. D. Moor, J. Vandewalle, A multilinear singular value decomposition, *SIAM Journal on Matrix Analysis & Applications* 21 (4) (2000) 1253–1278.
- [25] T. G. Kolda, B. W. Bader, Tensor decompositions and applications, *SIAM Review* 51 (3) (2009) 455–500.
- [26] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al., Distributed optimization and statistical learning via the alternating direction method of multipliers, *Foundations and Trends® in Machine learning* 3 (1) (2011) 1–122.
- [27] D. P. Bertsekas, *Nonlinear programming*, Athena scientific Belmont, 1999.
- [28] S. Arlot, A. Celisse, et al., A survey of cross-validation procedures for model selection, *Statistics surveys* 4 (2010) 40–79.

- [29] M. Steinlechner, Riemannian optimization for high-dimensional tensor completion, *SIAM Journal on Scientific Computing* 38 (5) (2016) 461–484.
- [30] W. Deng, W. Yin, On the global and linear convergence of the generalized alternating direction method of multipliers, *Journal of Scientific Computing* 66 (3) (2016) 889–916.
- [31] M. H. DeGroot, M. J. Schervish, *Probability and statistics*, Pearson Education, 2012.
- [32] J.-P. Ryckaert, G. Ciccotti, H. Berendsen, Numerical integration of the cartesian equations of motion of a system with constraints: molecular dynamics of n-alkanes, *Journal of Computational Physics* 23 (3) (1977) 327–341.
- [33] H. Elman, D. Silvester, A. Wathen, *Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics*, Oxford University Press (UK), 2014.
- [34] H. Elman, A. Ramage, D. Silvester, IFISS: A computational laboratory for investigating incompressible flow problems, *SIAM Review* 56 (2014) 261–273.