

EPFL

MASTER THESIS

Randomized low-rank approximation of matrices and tensors

Author:
Davide PRADOVERA

Supervisor:
Prof. Daniel KRESSNER

Examiner:
Dr. Michael PETERS

*A thesis submitted in fulfillment of the requirements
for the Master degree in Computational Science and Engineering*

within the

Numerical Analysis and High-Performance Computing Group

Section of Mathematics

School of Basic Sciences

Submitted on June 23, 2017 (v. July 11, 2017)



EPFL

Abstract

School of Basic Sciences
Section of Mathematics
Numerical Analysis and High-Performance Computing Group

Master degree

Randomized low-rank approximation of matrices and tensors

by Davide PRADOVERA

Low-rank matrix and tensor approximation is a central topic in computational linear algebra, with countless applications in scientific computing and data analysis. In this thesis, we explore the possibility to compute low-rank approximations of matrices and tensors by using randomized techniques. In particular, we provide an overview of the existing literature on the topic, and we derive several extensions to tackle two variations of the compression problem: the computation of low-rank approximations with fixed rank and with fixed precision.

This thesis shows that randomized approaches provide a very appealing alternative to deterministic techniques. In particular, in the case of matrix compression, we describe methods which have (approximately) the same complexity as the most efficient deterministic algorithms, but are more appealing because of their robustness.

Regarding tensors, this thesis focuses on the Tensor Train (TT) format, which represents a generalization of the low-rank decomposition for matrices. In particular, we derive a randomized approach which performs better than the *state-of-the-art* deterministic algorithm for the compression of tensors in the TT format, by exploiting the efficiency with which some operations can be carried out in this format.

Moreover, this thesis extends such procedures to the compression of Hadamard products of matrices and of tensors, which represent fundamental operations in several linear algebra algorithms. In both cases, we devise randomized techniques which outperform their deterministic counterparts. Also, we show that randomness leads to a diminished memory exploitation.

To corroborate all these claims, we provide extensive numerical experiments.

Keywords: Low-rank approximation, principal component analysis, tensor decomposition, higher-order singular value decomposition (HOSVD), Tensor Train format, Hadamard product, randomized algorithms.

Acknowledgements

This thesis would not have been possible without the support and the contributions of many amazing people.

First of all, I would like to thank my thesis supervisor Prof. Daniel Kressner for allowing me to work on this incredibly interesting topic and for the many suggestions, which were always on target. I have learned a lot during the development of this project, and not only in computational linear algebra.

I am also very grateful towards Dr. Michael Peters for his helpful suggestions and support, and for agreeing to serve as my examiner.

Moreover, I would like to express my gratitude to all my fellow MA students at EPFL, who endured these last years with me. In particular, I want to thank Ondine and Stefano for always being there when I needed to take a break, as well as Alberto, Aleksa, Antonio, Carlos, Marco, Niccolò, Nicolò, Przemysław, Riccardo, Santo and Stefano for the many lunches together. I have really enjoyed the time spent with you, and hope we will stay in touch.

Farther away, but not less important, I am extremely grateful for the support I have received from my family: my mother, my grandmother, my aunt and my sister. Your encouragement has always been a great source of motivation for me.

Lastly, thank you, Chiara, for having been able to put up with me throughout these years, for helping me in finding and following my way, and for giving me a place to go back to when I am stranded.

Contents

| | |
|--------------------------------------------------------------------------|------------|
| Abstract | iii |
| Acknowledgements | v |
| Contents | vii |
| 1 Introduction | 1 |
| 2 Preliminaries | 3 |
| 2.1 Notation and basic operations with matrices | 3 |
| 2.2 Notation and basic operations with tensors | 5 |
| 3 Matrix compression | 9 |
| 3.1 Randomized matrix compression with fixed rank | 10 |
| 3.1.1 Theoretical bounds for the randomized algorithm | 12 |
| 3.2 The inverse problem: low-rank approximation with fixed precision . . | 13 |
| 3.2.1 Bounds for the Frobenius norm | 14 |
| 3.2.2 Bounds for the spectral norm | 17 |
| 3.3 Numerical examples | 20 |
| 3.3.1 Computational environment | 20 |
| 3.3.2 Discretization of a rational function | 20 |
| 3.3.3 Compression of a Matérn kernel | 21 |
| 3.3.4 Compression with non-Gaussian random vectors | 23 |
| 4 Compression of tensors in the TT format | 25 |
| 4.1 Overview of the Tensor Train format | 25 |
| 4.2 Compression of tensors in the TT format using deterministic methods | 28 |
| 4.2.1 Orthogonalization of tensors in the TT format | 29 |
| 4.2.2 Truncation of tensors in the TT format in two sweeps | 31 |
| 4.2.3 Rounding of tensors in the TT format in two sweeps | 35 |
| 4.3 Compression of tensors in the TT format using a randomized approach | 37 |
| 4.3.1 Multidimensional tensor contraction | 37 |
| 4.3.2 Truncation of tensors in the TT format in a single sweep | 41 |
| 4.3.3 Randomized rounding of tensors in the TT format | 44 |
| 4.4 Numerical examples | 46 |
| 4.4.1 Truncation of a Scholes-like tensor | 46 |
| 4.4.2 Truncation of the discretization of a rational function | 48 |
| 4.4.3 Rounding of a random tensor | 49 |

| | | |
|----------|--------------------------------------------------------------------------------|-----------|
| 5 | Compression of Hadamard products | 51 |
| 5.1 | Compression of Hadamard products of matrices | 51 |
| 5.2 | Compression of Hadamard products of tensors in the TT format | 54 |
| 5.2.1 | Sampling the range of the unfoldings of Hadamard products . . | 55 |
| 5.2.2 | Low-rank approximation of tensors in the TT format | 56 |
| 5.3 | Numerical examples | 60 |
| 5.3.1 | Square of a 2-dimensional rational function through a Hadamard power | 60 |
| 5.3.2 | Hadamard product of high-dimensional rational functions . . . | 62 |
| 5.3.3 | Hadamard product of random tensors | 63 |
| 6 | Conclusion | 65 |
| | Bibliography | 67 |

Dedicated to my father

Chapter 1

Introduction

Matrices with low numerical rank are ubiquitous in a plethora of applications in mathematics, engineering and computer science, among which we cite just a few names: recommender systems, principal component analysis (PCA), compressive sampling, numerical solution of PDEs.

For such matrices, it is often appealing to derive a low-rank approximation of the form

$$\mathbf{A} \approx \mathbf{W}\mathbf{Z}^\top \tag{1.1}$$

with the two factors \mathbf{W} and \mathbf{Z} being “tall” matrices with just a few columns.

This representation is desirable for several reasons. Foremost, the compressed matrix can be stored very efficiently, and we can compute matrix-vector products quite cheaply. Moreover, the two factors provide an insight on the range and on the row space of the matrix: indeed, if the approximation (1.1) is sufficiently precise, the result obtained by applying \mathbf{A} (resp. \mathbf{A}^\top) onto a general vector lies mostly within the range of \mathbf{W} (resp. \mathbf{Z}).

Several decades of research have led to many classical methods which can be applied to obtain representation (1.1): these include the singular value decomposition (SVD), the interpolative decomposition (ID) and the rank-revealing QR (RRQR) decomposition. Each method comes with its own advantages and disadvantages: in particular, it has proven complicated to design algorithms which are both efficient and robust.

More recently, on the wake of a rising interest in randomized approaches to numerical problems, Halko et al. [HMT11] have presented a stochastic algorithm to tackle the low-rank approximation problem, which achieves both efficiency and robustness, at the cost of a (slightly) increased approximation error. In this thesis we perform a brief review of this randomized procedure, and generalize it to deal with different versions of the compression problem.

In particular, we derive a randomized approach to compute low-rank approximations of *tensors*. Tensors are the natural generalizations of matrices to dimensions higher than 2. They provide a very compact representation of high-dimensional data, which is useful in many applications in a wide variety of fields, for instance in the numerical approximation of stochastic PDEs (see e.g. [BNT07]) and in molecular quantum dynamics (see e.g. [Lub08] and [Ven+07]).

For tensors, a low-rank representation generalizing (1.1) is even more appealing, due to the countless issues which arise when increasing the number of dimensions (the order) of the tensor. Among these, the most critical is probably the *curse of*

dimensionality, which states that most operations involving tensors have a computational cost increasing exponentially in the order. At a glance, this may not seem too troublesome, but many applications involve tensors of order 100 or 1000: in these cases it is impossible even to store such tensors on most computers!

For this reason, in this thesis we mostly consider tensors in the Tensor Train (TT) format, which can be manipulated (using the majority of the basic operations) without being affected by the *curse of dimensionality*. For tensors in this format, we devise algorithms which outperform the *state-of-the-art* deterministic approach.

Moreover, we devise randomized techniques to compress Hadamard products of tensors, which represent one of the most basic operations in tensor-based computations. In this case we derive an algorithm whose memory requirements and complexity are considerably lower than the ones of its deterministic counterpart, as shown by several numerical experiments.

Chapter 2

Preliminaries

2.1 Notation and basic operations with matrices

We start by summarizing the notation employed for matrices and vectors. Throughout this thesis, we represent vectors and matrices with bold letters, respectively lower- and upper-case, e.g. $\mathbf{v} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$ for some positive integers m and n .

In particular, we indicate with \mathbf{I} the identity matrix, whose size we often clarify by using a subscript, e.g. $\mathbf{I}_m \in \mathbb{R}^{m \times m}$.

Elements of vectors and matrices are referred to in a MATLAB[®]-like format, e.g. $\mathbf{v}(i)$ and $\mathbf{A}(i, j)$ for suitable indices i and j . Moreover, we use a MATLAB[®] format to access sub-vectors and sub-matrices as well. For instance, if $\mathbf{A} \in \mathbb{R}^{20 \times 20}$, we have the following block representations:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}(1:5, 1:10) & \mathbf{A}(1:5, 11:20) \\ \mathbf{A}(6:20, 1:10) & \mathbf{A}(6:20, 11:20) \end{bmatrix} = \begin{bmatrix} \mathbf{A}(1:5, :) \\ \mathbf{A}(6:20, :) \end{bmatrix}$$

In the later chapters we also make use of the `diag` operation from MATLAB[®], which is defined as follows. Given a vector $\mathbf{v} \in \mathbb{R}^m$, we define the matrix $\text{diag}(\mathbf{v}) \in \mathbb{R}^{m \times m}$ as

$$(\text{diag}(\mathbf{v}))(i, j) = \mathbf{v}(i)\mathbf{I}_m(i, j) \quad \text{for } (i, j) \in \{1, \dots, m\}^2$$

Moreover, given a square matrix $\mathbf{A} \in \mathbb{R}^{m \times m}$, we define the vector $\text{diag}(\mathbf{A}) \in \mathbb{R}^m$ as

$$(\text{diag}(\mathbf{A}))(i) = \mathbf{A}(i, i) \quad \text{for } i \in \{1, \dots, m\}$$

Given a vector $\mathbf{v} \in \mathbb{R}^m$, we represent as $\|\mathbf{v}\|$ its 2-norm

$$\|\mathbf{v}\| = \left(\sum_{i=1}^m \mathbf{v}(i)^2 \right)^{1/2} \quad (2.1)$$

whereas for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we employ the notations $\|\mathbf{A}\|$ and $\|\mathbf{A}\|_F$ respectively for the spectral norm

$$\|\mathbf{A}\| = \max_{\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \quad (2.2)$$

and for the Frobenius norm

$$\|\mathbf{A}\|_F = \left(\sum_{j=1}^n \|\mathbf{A}(:, j)\|^2 \right)^{1/2} \quad (2.3)$$

It is trivial to verify that the 2-norm and both matrix norms are unitarily invariant: given a matrix U with orthonormal columns (i.e. such that $U^\top U = I$), it holds that

$$\|v\| = \|Uv\| \quad \text{and} \quad \|A\| = \|UA\| \quad \text{for any } v \text{ and } A \text{ of suitable size}$$

with $\|\cdot\|$ being either the Frobenius or the spectral norm. Moreover, we will make use of the following property.

Proposition 2.1 (Submultiplicativity). *Let $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times k}$. Then*

$$\|AB\| \leq \|A\| \|B\| \tag{2.4}$$

with $\|\cdot\|$ being either the Frobenius or the spectral norm.

Proof. Let us consider the spectral norm first.

If B has rank 0, $\|AB\| = \|B\| = 0$ and the claim holds trivially. Otherwise, we may restrict the maximum in definition (2.2) to vectors which do not belong to the kernel of B , to obtain

$$\begin{aligned} \|AB\| &= \max_{\substack{x \in \mathbb{R}^k \setminus \{0\} \\ Bx \neq 0}} \frac{\|ABx\|}{\|x\|} \\ &\leq \max_{\substack{x \in \mathbb{R}^k \setminus \{0\} \\ Bx \neq 0}} \frac{\|ABx\|}{\|Bx\|} \max_{\substack{x \in \mathbb{R}^k \setminus \{0\} \\ Bx \neq 0}} \frac{\|Bx\|}{\|x\|} \\ &\leq \max_{y \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ay\|}{\|y\|} \max_{x \in \mathbb{R}^k \setminus \{0\}} \frac{\|Bx\|}{\|x\|} = \|A\| \|B\| \end{aligned}$$

If $\|\cdot\|$ is the Frobenius norm, we can apply the submultiplicativity of the spectral norm (or equivalently its definition) to obtain

$$\|AB\|_F^2 = \sum_{j=1}^k \|AB(:,j)\|^2 \leq \sum_{j=1}^k \|A\|^2 \|B(:,j)\|^2 = \|A\|^2 \|B\|_F^2$$

□

We also introduce several kinds of matrix-matrix multiplication. Consider $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{m' \times n'}$. We define:

- the Kronecker product $A \otimes B$ as the $(mm') \times (nn')$ matrix with entries

$$(A \otimes B)(i' + m'(i-1), j' + n'(j-1)) = A(i, j)B(i', j') \tag{2.5}$$

- if $m = m'$ and $n = n'$, the Hadamard product $A * B$ as the $m \times n$ matrix with entries

$$(A * B)(i, j) = A(i, j)B(i, j) \tag{2.6}$$

- if $n = n'$, the Khatri-Rao product $A \odot B$ as the $(mm') \times n$ matrix with entries

$$(A \odot B)(i' + m'(i-1), j) = A(i, j)B(i', j) \tag{2.7}$$

- if $m = m'$, the transpose Khatri-Rao product $\mathbf{A} \odot^\top \mathbf{B}$ as the $m \times (nm')$ matrix given by

$$\mathbf{A} \odot^\top \mathbf{B} = (\mathbf{A}^\top \odot \mathbf{B}^\top)^\top \quad (2.8)$$

In all the definitions above, the indices i, j, i' and j' range over all the possible values for which $\mathbf{A}(i, j)$ and $\mathbf{B}(i', j')$ make sense.

In the following chapters, we will make extensive use of the bilinearity of all the products defined above, as well as of the following properties of the Kronecker product, which can be proven quite easily by comparing the entries of the two sides of the equations.

Proposition 2.2. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m' \times n'}$, $\mathbf{C} \in \mathbb{R}^{n \times k}$ and $\mathbf{D} \in \mathbb{R}^{n' \times k'}$. Then*

$$(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top \quad (2.9)$$

and

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D}) \quad (2.10)$$

Also, it holds that

$$\mathbf{I}_m \otimes \mathbf{I}_n = \mathbf{I}_{mn} \quad (2.11)$$

2.2 Notation and basic operations with tensors

In this section we review the basic notation for tensors. For any integer $d > 0$ and any tuple of positive integers (n_1, \dots, n_d) of size d , we represent with a calligraphic bold upper-case letter the tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$. We indicate single entries of the tensor by generalizing the notation used for matrices: for a multi-index $(i_1, \dots, i_d) \in \{1, \dots, n_1\} \times \dots \times \{1, \dots, n_d\}$, the corresponding entry of the tensor is $\mathcal{A}(i_1, \dots, i_d)$.

The extraction of rows and columns can be generalized to tensors, again by using a MATLAB[®]-like notation: given a d -dimensional tensor \mathcal{A} , we may fix any number of indices to obtain a tensor with reduced order, by collapsing some of the modes. For instance, given $\mathcal{A} \in \mathbb{R}^{3 \times 2 \times 4 \times 6}$, we may consider the matrix $\mathcal{A}(:, 1, :, 4) \in \mathbb{R}^{3 \times 4}$, given by

$$\mathcal{A}(:, 1, :, 4) = \begin{bmatrix} \mathcal{A}(1, 1, 1, 4) & \mathcal{A}(1, 1, 2, 4) & \mathcal{A}(1, 1, 3, 4) & \mathcal{A}(1, 1, 4, 4) \\ \mathcal{A}(2, 1, 1, 4) & \mathcal{A}(2, 1, 2, 4) & \mathcal{A}(2, 1, 3, 4) & \mathcal{A}(2, 1, 4, 4) \\ \mathcal{A}(3, 1, 1, 4) & \mathcal{A}(3, 1, 2, 4) & \mathcal{A}(3, 1, 3, 4) & \mathcal{A}(3, 1, 4, 4) \end{bmatrix}$$

If we fix all the indices except one (resp. two), we talk about *fibers* (resp. *slices*) of the tensor.

Moreover, tensors can be reshaped into matrices using the following operation.

Definition 2.1 (Tensor unfolding). *Given a d -dimensional tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $j \in \{1, \dots, d-1\}$, the j^{th} unfolding of \mathcal{A} is the matrix $\mathcal{A}^{<j>} \in \mathbb{R}^{(n_1 \dots n_j) \times (n_{j+1} \dots n_d)}$ whose entries are*

$$\mathcal{A}^{<j>}(i_{\text{row}}(i_1, \dots, i_j), i_{\text{col}}(i_{j+1}, \dots, i_d)) = \mathcal{A}(i_1, \dots, i_d) \quad (2.12)$$

for $(i_1, \dots, i_d) \in I_1 \times \dots \times I_d$. The maps

$$i_{\text{row}} : I_1 \times \dots \times I_j \rightarrow \{1, \dots, n_1 \dots n_j\} \text{ and } i_{\text{col}} : I_{j+1} \times \dots \times I_d \rightarrow \{1, \dots, n_{j+1} \dots n_d\}$$

define a colexicographical ordering on $I_1 \times \dots \times I_j$ and $I_{j+1} \times \dots \times I_d$ respectively, i.e.

$$i_{\text{row}}(i_1, \dots, i_j) = 1 + \sum_{\mu=1}^j (i_\mu - 1) \prod_{\nu=1}^{\mu-1} n_\nu$$

$$i_{\text{col}}(i_{j+1}, \dots, i_d) = 1 + \sum_{\mu=j+1}^d (i_\mu - 1) \prod_{\nu=j+1}^{\mu-1} n_\nu$$

We additionally define the degenerate 0th unfolding $\mathcal{A}^{<0>}$ as the 1st unfolding of the tensor $\mathcal{A}' \in \mathbb{R}^{1 \times n_1 \times \dots \times n_d}$, defined entry-wise as

$$\mathcal{A}'(1, i_1, \dots, i_d) = \mathcal{A}(i_1, \dots, i_d) \quad \text{for } (i_1, \dots, i_d) \in I_1 \times \dots \times I_d$$

and the d^{th} unfolding $\mathcal{A}^{<d>}$ as the d^{th} unfolding of the tensor $\mathcal{A}'' \in \mathbb{R}^{n_1 \times \dots \times n_d \times 1}$, defined entry-wise as

$$\mathcal{A}''(i_1, \dots, i_d, 1) = \mathcal{A}(i_1, \dots, i_d) \quad \text{for } (i_1, \dots, i_d) \in I_1 \times \dots \times I_d$$

In particular, $\mathcal{A}^{<0>} \in \mathbb{R}^{1 \times (n_1 \dots n_d)}$ and $\mathcal{A}^{<d>} \in \mathbb{R}^{(n_1 \dots n_d) \times 1}$.

It is trivial to see that, if we fix the size of a d -dimensional tensor \mathcal{A} and an index $j \in \{0, \dots, d\}$, the j^{th} unfolding $\mathcal{A}^{<j>}$ uniquely identifies \mathcal{A} . Equivalently, for any $j \in \{0, \dots, d\}$, the “ j^{th} unfolding” operation defines a bijection from $\mathbb{R}^{n_1 \times \dots \times n_d}$ to $\mathbb{R}^{(n_1 \dots n_j) \times (n_{j+1} \dots n_d)}$.

In MATLAB[®] (and in any environment which relies on column-major array storage), the computation of any unfolding of a full tensor can be carried out without moving entries in memory, since data is stored in colexicographical order.

In particular, the conversion from tensor to matrix format (and vice versa) can be easily carried out using the `reshape` command. For instance, the j^{th} unfolding (for $j \in \{0, \dots, d\}$) of a dimensional tensor \mathcal{A} of size $\mathbf{n} = (n_1, \dots, n_d)$ can be obtained as

$$\text{reshape}(\mathcal{A}, [\text{prod}(\mathbf{n}(1:j)), \text{prod}(\mathbf{n}(j+1:d))])$$

In addition, for any tensor \mathcal{A} we define the Frobenius norm

$$\|\mathcal{A}\|_F = \left(\sum_{i_1=1}^{n_1} \dots \sum_{i_d=1}^{n_d} \mathcal{A}(i_1, \dots, i_d)^2 \right)^{1/2} \quad (2.13)$$

and we generalize the matrix-vector product operation as follows.

Definition 2.2 (*n*-mode product). Given a d -dimensional tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$, an integer $j \in \{1, \dots, d\}$ and a matrix $\mathbf{M} \in \mathbb{R}^{m \times n_j}$, the j^{th} mode product between \mathcal{A} and \mathbf{M} is the d -dimensional tensor $(\mathcal{A} \times_j \mathbf{M}) \in \mathbb{R}^{n_1 \times \dots \times n_{j-1} \times m \times n_{j+1} \times \dots \times n_d}$ such that

$$(\mathcal{A} \times_j \mathbf{M})(i_1, \dots, i_{j-1}, \alpha, i_{j+1}, \dots, i_d) = \sum_{i_j=1}^{n_j} \mathcal{A}(i_1, \dots, i_d) \mathbf{M}(\alpha, i_j) \quad (2.14)$$

To conclude, we observe that the Kronecker and Hadamard products can be naturally generalized to tensors. Given $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $\mathcal{B} \in \mathbb{R}^{n'_1 \times \dots \times n'_d}$, their

Kronecker product is the tensor $\mathcal{A} \otimes \mathcal{B} \in \mathbb{R}^{(n_1 n'_1) \times \dots \times (n_d n'_d)}$ with entries

$$(\mathcal{A} \otimes \mathcal{B})(i'_1 + n'_1(i_1 - 1), \dots, i'_d + n'_d(i_d - 1)) = \mathcal{A}(i_1, \dots, i_d) \mathcal{B}(i'_1, \dots, i'_d) \quad (2.15)$$

for all suitable multi-indices (i_1, \dots, i_d) and (i'_1, \dots, i'_d) .

Moreover, if \mathcal{A} and \mathcal{B} have the same size, then we can define their Hadamard product $\mathcal{A} * \mathcal{B} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ entry-wise as

$$(\mathcal{A} * \mathcal{B})(i_1, \dots, i_d) = \mathcal{A}(i_1, \dots, i_d) \mathcal{B}(i_1, \dots, i_d) \quad (2.16)$$

Chapter 3

Matrix compression

The general low-rank approximation problem for matrices may be formulated as follows: given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ (we assume $m \geq n$ without loss of generality) and an integer $r \leq n$, we look for two matrices $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, such that the error

$$\epsilon(\mathbf{W}, \mathbf{Z}) = \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \quad (3.1)$$

is minimized, with $\|\cdot\|$ being either the spectral or the Frobenius norm.

For unitarily invariant norms (e.g. the two norms considered above), the Eckart-Young and Mirsky theorems (see e.g. [EY36] and [Mir60]) guarantee that the problem of finding the two factors \mathbf{W} and \mathbf{Z} always admits an optimal solution. In particular, one such solution may be found by truncating (at rank r) the singular value decomposition of \mathbf{A} : indeed, if

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \quad (3.2)$$

is the SVD of \mathbf{A} , then the minimum of ϵ over $\mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$ is equal to

$$\epsilon(\mathbf{W}_{(r)}, \mathbf{Z}_{(r)}) = \|\mathbf{\Sigma}((r+1) : n, (r+1) : n)\|$$

with the optimal factors being

$$\mathbf{W}_{(r)} = \mathbf{U}(:, 1 : r) \quad \text{and} \quad \mathbf{Z}_{(r)} = \mathbf{V}(:, 1 : r)\mathbf{\Sigma}(1 : r, 1 : r) \quad (3.3)$$

It is trivial to see that this solution is not unique (except when \mathbf{A} has rank 0). Indeed, for any non-singular matrix $\mathbf{M} \in \mathbb{R}^{r \times r}$, the factors

$$\mathbf{W}_{(r)}^* = \mathbf{W}_{(r)}\mathbf{M} \quad \text{and} \quad \mathbf{Z}_{(r)}^* = \mathbf{Z}_{(r)}\mathbf{M}^{-\top}$$

are still optimal, since $\mathbf{W}_{(r)}^* \mathbf{Z}_{(r)}^{*\top} = \mathbf{W}_{(r)} \mathbf{Z}_{(r)}^\top$.

Given this result, a naive algorithm to find a low-rank approximation for a general unstructured matrix of size $m \times n$ may be based on the explicit computation (and successive truncation) of the full SVD of the matrix. The resulting computational cost is $\mathcal{O}(mn^2)$ flops, independently of r .

More advanced techniques rely on a RRQR factorization (see e.g. [Cha87]) to identify the dominant range of the matrix, and then manipulate the results to obtain the desired low-rank approximation. For most matrices these approaches are quite efficient, having a complexity of $\mathcal{O}(mnr)$ flops, but in some extreme cases (see e.g. [GE96]) their cost is the same as the computation of the full SVD, i.e. $\mathcal{O}(mn^2)$ flops.

In comparison, the randomized algorithm presented in the upcoming section has a computational cost approximately equal to $\mathcal{O}(mnr)$, and has the benefit of being robust. The disadvantage of this stochastic approach is, of course, the presence of an intrinsically positive probability that the algorithm will fail.

3.1 Randomized matrix compression with fixed rank

Due to the orthonormality of the columns of U , it is trivial to see that the optimal low-rank approximation (3.3) of the matrix A may be obtained by projecting the matrix onto the range of $U(:, 1:r)$, i.e. the span of the r dominant left-singular vectors of A :

$$\mathbf{W}_{(r)}\mathbf{Z}_{(r)}^\top = U(:, 1:r)\Sigma(1:r, 1:r)V(:, 1:r)^\top = U(:, 1:r)U(:, 1:r)^\top A$$

Hence, we can find a low-rank approximation in two steps: first, we compute a left factor W (with full rank) whose range is an approximation of the range of $U(:, 1:r)$. Then we find the corresponding right factor Z as

$$Z = A^\top W (W^\top W)^{-1}$$

so that

$$WZ^\top = W(W^\top W)^{-1}W^\top A \approx U(:, 1:r)U(:, 1:r)^\top A$$

since $W(W^\top W)^{-1}W^\top$ represents the projection matrix onto the range of W .

In particular, we remark that the term $(W^\top W)^{-1}$ may be omitted if W has orthonormal columns. To avoid the computation of this term, in the following we will always require that the left factor W is properly orthogonalized.

Thus, the problem is now reduced to finding an approximation of the dominant subspace (of dimension r) of the range of A . To this aim, we follow the intuition provided in [HMT11]: let

$$E = A - \mathbf{W}_{(r)}\mathbf{Z}_{(r)}^\top$$

represent the perturbation from the desired low-rank approximation. In particular, due to the properties of the truncated SVD, the range of E and the range of $W_{(r)}$ are orthogonal, i.e. $W_{(r)}^\top E = 0$.

Given a family of independent Gaussian random vectors $\{\omega^{(l)}\}_{l=1}^n \subset \mathbb{R}^n$, we may randomly sample the range of A as follows:

$$\mathbf{y}^{(l)} = A\omega^{(l)} = \mathbf{W}_{(r)}\mathbf{Z}_{(r)}^\top\omega^{(l)} + E\omega^{(l)} \quad \text{for } l = 1, \dots, n$$

Almost surely, the action of the perturbation E onto the random vectors makes each sample fall outside the range of $W_{(r)}$. Accordingly, depending on the relative magnitude of the singular values of A , the span of just r samples $\{\mathbf{y}^{(l)}\}_{l=1}^r$ may provide quite a poor approximation for the range of $W_{(r)}$. To reduce the error, we can set an oversampling parameter $p > 0$ and use an enriched set of samples $\{\mathbf{y}^{(l)}\}_{l=1}^{r+p}$, which yields a much better approximation (see Theorem 3.2).

Now we compute (e.g. through a QR decomposition) the matrix $\widetilde{W} \in \mathbb{R}^{m \times (r+p)}$, whose columns form an orthonormal basis of the span of $\{\mathbf{y}^{(l)}\}_{j=1}^{r+p}$. This leads to a

rank- $(r + p)$ approximation of the form

$$\mathbf{A} \approx \tilde{\mathbf{A}} = \tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top \mathbf{A} = \tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top$$

In order to shrink this approximation to rank r , we carry out a *recompression* step by computing the best rank- r approximation of $\tilde{\mathbf{A}}$. In particular, we observe that, since $\tilde{\mathbf{W}}$ has orthonormal columns, the SVD of $\tilde{\mathbf{A}}$ is of the form

$$\tilde{\mathbf{A}} = (\tilde{\mathbf{W}}\tilde{\mathbf{U}})\tilde{\Sigma}\tilde{\mathbf{V}}^\top$$

for some orthogonal matrix $\tilde{\mathbf{U}} \in \mathbb{R}^{(r+p) \times (r+p)}$. Hence, it suffices to find the best rank- r approximation of $\tilde{\mathbf{Z}}^\top$, and then multiply it by $\tilde{\mathbf{W}}$ from the left.

This time we can explicitly compute the SVD of the “tall” matrix $\tilde{\mathbf{Z}} \in \mathbb{R}^{n \times (r+p)}$, and truncate it at rank r . The resulting factors for the approximation of \mathbf{A} are

$$\mathbf{W} = \tilde{\mathbf{W}}\tilde{\mathbf{U}}(:, 1:r) \quad \text{and} \quad \mathbf{Z} = \tilde{\mathbf{V}}(:, 1:r)\tilde{\Sigma}(1:r, 1:r)$$

with $\tilde{\mathbf{V}}\tilde{\Sigma}\tilde{\mathbf{U}}^\top$ being the SVD of $\tilde{\mathbf{Z}}$.

The resulting procedure is summarized in Algorithm 1.

Algorithm 1: Randomized compression with fixed rank

Data: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, positive integers r and p (with $r < n$ and $p + r \leq n$)

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$ and $\mathbf{A} \approx \mathbf{W}\mathbf{Z}^\top$

- 1 Generate Gaussian matrix $\Omega \in \mathbb{R}^{n \times (r+p)}$
 - 2 $\mathbf{Y} := \mathbf{A}\Omega$
 - 3 Compute QR decomposition: $\mathbf{Y} =: \tilde{\mathbf{W}}\tilde{\mathbf{R}}$
 - 4 $\tilde{\mathbf{Z}} := \mathbf{A}^\top \tilde{\mathbf{W}}$
 - 5 Compute SVD: $\tilde{\mathbf{Z}} =: \tilde{\mathbf{V}}\tilde{\Sigma}\tilde{\mathbf{U}}^\top$
 - 6 $\mathbf{W} := \tilde{\mathbf{W}}\tilde{\mathbf{U}}(:, 1:r)$
 - 7 $\mathbf{Z} := \tilde{\mathbf{V}}(:, 1:r)\tilde{\Sigma}(1:r, 1:r)$
-

Assuming that \mathbf{A} is treated as a general dense matrix, the total computational cost is

$$\mathcal{O}(n(r+p)T_{rng} + mn(r+p))$$

flops, with T_{rng} being the cost of drawing a single random sample. If we assume (as we will do from here onward) that the cost of generating random samples is negligible, the overall complexity is $\mathcal{O}(mn(r+p))$. In particular, when $r + p \ll n$ (which is usually the case in practice), the complexity of the algorithm is dominated by the cost of computing the two products $\mathbf{A}\Omega$ and $\mathbf{A}^\top \tilde{\mathbf{W}}$.

More generally, if the costs of matrix-vector multiplications by \mathbf{A} and \mathbf{A}^\top are respectively T_{mult} and T'_{mult} , the complexity of the algorithm is

$$\mathcal{O}((T_{mult} + T'_{mult})(r+p) + m(r+p)^2)$$

3.1.1 Theoretical bounds for the randomized algorithm

The compressed matrix resulting from Algorithm 1 satisfies the following *a priori* bounds.

Theorem 3.1. *Let \mathbf{W} and \mathbf{Z} be the two factors computed by Algorithm 1, and assume $p \geq 2$. If $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ are the singular values of \mathbf{A} , then*

$$\mathbb{E} \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|_F \leq \left(1 + \sqrt{1 + \frac{r}{p-1}}\right) \left(\sum_{j=r+1}^n \sigma_j^2\right)^{1/2} \quad (3.4)$$

and

$$\mathbb{E} \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \leq \left(2 + \sqrt{\frac{r}{p-1}}\right) \sigma_{r+1} + \frac{e\sqrt{r+p}}{p} \left(\sum_{j=r+1}^n \sigma_j^2\right)^{1/2} \quad (3.5)$$

Proof. Let $\tilde{\mathbf{W}}$ and $\tilde{\mathbf{Z}}$ be defined as in Algorithm 1.

Using the triangular inequality we can write

$$\begin{aligned} \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| &\leq \|\mathbf{A} - \tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top\| + \|\tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top - \mathbf{W}\mathbf{Z}^\top\| \\ &= \|\mathbf{A} - \tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top\| + \|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \end{aligned}$$

with $\|\cdot\|$ being either the Frobenius or spectral norm.

We can bound the expected value of the first term by using Theorems 10.5 and 10.6 in [HMT11], which state that

$$\mathbb{E} \|\mathbf{A} - \tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top\|_F \leq \left(1 + \frac{r}{p-1}\right)^{1/2} \left(\sum_{j=r+1}^n \sigma_j^2\right)^{1/2}$$

and

$$\mathbb{E} \|\mathbf{A} - \tilde{\mathbf{W}}\tilde{\mathbf{Z}}^\top\| \leq \left(1 + \sqrt{\frac{r}{p-1}}\right) \sigma_{r+1} + \frac{e\sqrt{r+p}}{p} \left(\sum_{j=r+1}^n \sigma_j^2\right)^{1/2}$$

To bound the second term, we observe that $\mathbf{W}\mathbf{Z}^\top$ is, by construction, the best rank- r approximation of $\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A}$ (in both the Frobenius and the spectral norm, due to the Mirsky theorem). Hence

$$\|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| = \min_{\substack{\mathbf{B} \in \mathbb{R}^{m \times n} \\ \text{rank } \mathbf{B} \leq r}} \|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A} - \mathbf{B}\| \leq \|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A} - \mathbf{C}\|$$

for any $\mathbf{C} \in \mathbb{R}^{m \times n}$ with rank at most r .

In particular, we may pick $\mathbf{C} = \tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A}_{(r)}$, with $\mathbf{A}_{(r)}$ being the best rank- r approximation of \mathbf{A} . This yields

$$\|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \leq \|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top(\mathbf{A} - \mathbf{A}_{(r)})\|$$

By orthogonality, it holds that

$$\|\tilde{\mathbf{W}}\tilde{\mathbf{W}}^\top\| = \|\mathbf{I}_{r+p}\| = 1$$

and we can apply the submultiplicativity of $\|\cdot\|$ (see Proposition 2.1) to obtain

$$\|\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^\top (\mathbf{A} - \mathbf{A}_{(r)})\| \leq \|\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^\top\| \|\mathbf{A} - \mathbf{A}_{(r)}\| = \|\mathbf{A} - \mathbf{A}_{(r)}\|$$

or equivalently

$$\|\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^\top (\mathbf{A} - \mathbf{A}_{(r)})\|_F \leq \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2} \quad \text{and} \quad \|\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^\top (\mathbf{A} - \mathbf{A}_{(r)})\| \leq \sigma_{r+1}$$

Hence, the expected value of the residual can be bounded as

$$\mathbb{E}\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|_F \leq \left(1 + \frac{r}{p-1}\right)^{1/2} \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2} + \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2}$$

and

$$\mathbb{E}\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \leq \left(1 + \sqrt{\frac{r}{p-1}}\right) \sigma_{r+1} + \frac{e\sqrt{r+p}}{p} \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2} + \sigma_{r+1}$$

□

For fixed values of r and p , Theorem 3.1 proves that Algorithm 1 is (in average) quasi-optimal in the Frobenius norm, but not in the spectral norm. For this reason, the application of the Frobenius norm appears more natural in the framework of randomized low-rank approximation.

We also have the following estimates on the failure probability of the algorithm.

Theorem 3.2. *Consider the hypotheses of Theorem 3.1, with $p \geq 4$. Then it holds*

$$\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|_F \leq \left(2 + 10\sqrt{\frac{r}{p}}\right) \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2} + 11\sqrt{1 + \frac{r}{p}} \sigma_{r+1} \quad (3.6)$$

except with probability $7e^{-p}$ at most, and

$$\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| \leq 8\frac{\sqrt{r+p}}{p+1} \left(\sum_{j=r+1}^n \sigma_j^2 \right)^{1/2} + \left(2 + 20\sqrt{1 + \frac{r}{p}}\right) \sigma_{r+1} \quad (3.7)$$

except with probability $6e^{-p}$ at most.

Proof. The proof is analogous to the one of Theorem 3.1, with the only difference that we use Theorems 10.7 and 10.8 (with the values $t = e$ and $u = \sqrt{2p}$) from [HMT11] instead of 10.5 and 10.6. □

3.2 The inverse problem: low-rank approximation with fixed precision

In most applications, the target rank r is not known in advance. Instead, we may wish to compress a matrix with fixed precision: given $\mathbf{A} \in \mathbb{R}^{m \times n}$ (with $m \geq n$) and

$\varepsilon > 0$, we seek $\mathbf{W} \in \mathbb{R}^{m \times \bar{r}}$ and $\mathbf{Z} \in \mathbb{R}^{n \times \bar{r}}$ such that

$$\frac{\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|}{\|\mathbf{A}\|} \leq \varepsilon \quad (3.8)$$

with \bar{r} being the smallest integer such that (3.8) is feasible.

A simple, deterministic way to solve this problem starts with the explicit computation of the SVD of \mathbf{A} . Then the value of \bar{r} can be deduced from its singular values ($\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$), as the smallest integer which satisfies

$$\sum_{j=\bar{r}+1}^n \sigma_j^2 \leq \varepsilon^2 \sum_{j=1}^n \sigma_j^2$$

for the Frobenius norm, or

$$\sigma_{\bar{r}+1} \leq \varepsilon \sigma_1$$

for the spectral norm, and the two factors can be computed by truncating the SVD at rank \bar{r} .

Similarly to the fixed-rank case, this process can be sped up by applying a RRQR decomposition. This typically reduces the complexity from $\mathcal{O}(mn^2)$ to $\mathcal{O}(mn\bar{r})$, at the cost of a decrease in the robustness of the algorithm.

A randomized approach, obtained as a suitable generalization of Algorithm 1, can achieve both efficiency and robustness. To derive such an algorithm, the main issue is finding an inexpensive way to compute or estimate the relative error (3.8). Indeed, once such an estimate is available, it suffices to apply the following adaptive approach: we proceed by using the framework described in Section 3.1, and we keep adding new samples from the range of \mathbf{A} until the approximation $\mathbf{W}\mathbf{Z}^\top$ achieves the desired precision.

Accordingly, we devote the upcoming sections to the description of several approaches to tackle this issue.

3.2.1 Bounds for the Frobenius norm

Assume that the relative precision is specified in terms of the Frobenius norm.

Let $\mathbf{W}\mathbf{Z}^\top$ be a rank- r approximation of \mathbf{A} , with \mathbf{W} having orthonormal columns, and let $\mathbf{Z} = \mathbf{A}^\top \mathbf{W}$. We wish to verify if (3.8) is satisfied.

To this aim, we can observe that the following orthogonal decomposition holds:

$$\mathbf{A} = \mathbf{A} - \mathbf{W}\mathbf{Z}^\top + \mathbf{W}\mathbf{Z}^\top = (\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A} + \mathbf{W}\mathbf{W}^\top \mathbf{A}$$

Indeed, $\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top$ and $\mathbf{W}\mathbf{W}^\top$ are projection matrices onto orthogonal subspaces of \mathbb{R}^m .

Thus, the Pythagorean theorem yields

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|_F^2 + \|\mathbf{W}\mathbf{Z}^\top\|_F^2 \quad (3.9)$$

and (3.8) is equivalent to

$$\|\mathbf{A}\|_F^2 - \|\mathbf{W}\mathbf{Z}^\top\|_F^2 \leq \varepsilon^2 \|\mathbf{A}\|_F^2 \quad (3.10)$$

Let us consider the term $\|\mathbf{W}\mathbf{Z}^\top\|_F$. Since \mathbf{W} has orthonormal columns, it is sufficient to compute the norm of the “tall” right factor \mathbf{Z} . If we also allow the explicit computation of $\|\mathbf{A}\|_F$ (which needs to be carried out only once), then we obtain the procedure summarized in Algorithm 2.

Algorithm 2: Randomized compression with fixed (Frobenius) precision with explicit computation of $\|\mathbf{A}\|_F$

Data: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\varepsilon > 0$

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ and such that (3.8) holds for the Frobenius norm

```

1 Compute  $\|\mathbf{A}\|_F$ 
2  $\mathbf{W} := []$ ,  $\mathbf{Z} := []$ 
3 while  $\|\mathbf{A}\|_F^2 - \|\mathbf{Z}\|_F^2 > \varepsilon^2 \|\mathbf{A}\|_F^2$  do
4   Generate Gaussian vector  $\boldsymbol{\omega} \in \mathbb{R}^n$ 
5    $\mathbf{w} := (\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A}\boldsymbol{\omega}$ 
6    $\mathbf{w} := \mathbf{w}/\|\mathbf{w}\|$ 
7    $\mathbf{W} := [\mathbf{W} \mid \mathbf{w}]$ 
8    $\mathbf{z} := \mathbf{A}^\top \mathbf{w}$ 
9    $\mathbf{Z} := [\mathbf{Z} \mid \mathbf{z}]$ 
10 end

```

Since each iteration adds a single column to \mathbf{Z} , the computation of $\|\mathbf{Z}\|_F$ can be carried out recursively in merely $\mathcal{O}(nr)$ total operations. If \mathbf{A} is stored as a full matrix, the computation of its Frobenius norm takes $\mathcal{O}(mn)$ flops, not affecting the overall complexity $\mathcal{O}(mnr)$, with the dominant operations being, as in the fixed-rank case, the multiplications $\mathbf{A}\boldsymbol{\omega}$ and $\mathbf{A}^\top \mathbf{w}$.

However, due to the randomness of the approach, with high probability the final rank r will be larger than the optimal one \bar{r} . In some cases (see e.g. Section 3.3.3), the algorithm may overshoot the optimal rank \bar{r} by a considerable amount. As such, it may be advisable to add a recompression step at the end of the algorithm:

1. Compute the SVD of \mathbf{Z} : $\mathbf{Z} = \mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^\top$.
2. Find the smallest r' which satisfies

$$\|\mathbf{A}\|_F^2 - \sum_{j=1}^{r'} \sigma_j^2 \leq \varepsilon^2 \|\mathbf{A}\|_F^2 \quad (3.11)$$

with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$ being the diagonal entries of $\boldsymbol{\Sigma}$.

3. Update $\mathbf{W}^* = \mathbf{W}\mathbf{U}(:, 1:r')$ and $\mathbf{Z}^* = \mathbf{V}(:, 1:r')\boldsymbol{\Sigma}(1:r', 1:r')$.

Error estimation without the norm of the full matrix

In many cases of interest, computing $\|\mathbf{A}\|_F$ is much more expensive than $\mathcal{O}(mn)$ flops. For instance, assume that \mathbf{A} is not stored in memory, and that we can access it only through matrix-vector multiplications. In this case, computing the Frobenius norm, or even estimating it with reasonable accuracy, would be much more expensive than the rest of the algorithm.

To tackle these situations, we proceed to reformulate (3.10) in a cheaper way. We start by providing a lower-bound for the denominator: using (3.9), we obtain

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\|_F^2 + \|\mathbf{W}\mathbf{Z}^\top\|_F^2 \geq \|\mathbf{W}\mathbf{Z}^\top\|_F^2 = \|\mathbf{Z}\|_F^2$$

In particular, as the algorithm goes on, the ranks of \mathbf{W} and \mathbf{Z} increase, and the discrepancy between $\|\mathbf{Z}\|_F$ and $\|\mathbf{A}\|_F$ becomes progressively smaller.

This fact proves extremely useful in the estimation of the numerator. Assume that $\widetilde{\mathbf{W}}$ is the matrix (with orthonormal columns) obtained by carrying out $q > 0$ more steps of the algorithm (i.e. it has q more columns than \mathbf{W} , and $\widetilde{\mathbf{W}} = [\mathbf{W} | \dots]$), and let $\widetilde{\mathbf{Z}} = \mathbf{A}^\top \widetilde{\mathbf{W}}$. Since $\widetilde{\mathbf{W}}\widetilde{\mathbf{Z}}^\top$ is a “richer” approximation of \mathbf{A} , we use its norm to approximate the value of $\|\mathbf{A}\|_F$.

In particular, a suitable value of q should be chosen using the prior knowledge on the spectrum of \mathbf{A} . If we expect the singular values of \mathbf{A} to decay exponentially (see e.g. Section 3.3.2), it is enough to choose a small q , whereas a larger value should be used for slower decays (see e.g. Section 3.3.3).

In summary, as an approximation of (3.10), we enforce the condition

$$\|\widetilde{\mathbf{Z}}\|_F^2 - \|\mathbf{Z}\|_F^2 \leq \varepsilon^2 \|\widetilde{\mathbf{Z}}\|_F^2 \quad (3.12)$$

This leads to Algorithm 3, whose complexity is again $\mathcal{O}(mnr)$.

Algorithm 3: Randomized compression with fixed (Frobenius) precision without computing $\|\mathbf{A}\|_F$

Data: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\varepsilon > 0$, integer $q > 0$

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ and such that (3.8) holds with high probability for the Frobenius norm

```

1 Generate Gaussian matrix  $\mathbf{\Omega} \in \mathbb{R}^{n \times q}$ 
2  $\mathbf{Y} := \mathbf{A}\mathbf{\Omega}$ 
3 Compute QR decomposition:  $\mathbf{Y} =: \mathbf{W}\mathbf{R}$ 
4  $\mathbf{Z} := \mathbf{A}^\top \mathbf{W}$ 
5  $r := q$ 
6 while  $\|\mathbf{Z}\|_F^2 - \|\mathbf{Z}(:, 1 : (r - q))\|_F^2 > \varepsilon^2 \|\mathbf{Z}\|_F^2$  do
7    $r := r + 1$ 
8   Generate Gaussian vector  $\boldsymbol{\omega} \in \mathbb{R}^n$ 
9    $\mathbf{w} := (\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A}\boldsymbol{\omega}$ 
10   $\mathbf{w} := \mathbf{w} / \|\mathbf{w}\|$ 
11   $\mathbf{W} := [\mathbf{W} | \mathbf{w}]$ 
12   $\mathbf{z} := \mathbf{A}^\top \mathbf{w}$ 
13   $\mathbf{Z} := [\mathbf{Z} | \mathbf{z}]$ 
14 end

```

In particular, we remark that $\widetilde{\mathbf{W}}$ and $\widetilde{\mathbf{Z}}$ are used instead of \mathbf{W} and \mathbf{Z} as the final low-rank factors, since $\widetilde{\mathbf{W}}\widetilde{\mathbf{Z}}^\top$ always represents a better approximation of \mathbf{A} .

As above, we may wish to add a recompression step at the end of the procedure. However, since this time the Frobenius norm of \mathbf{A} is not known, we need to replace $\|\mathbf{A}\|_F$ in (3.11) with $\|\widetilde{\mathbf{Z}}\|_F$.

3.2.2 Bounds for the spectral norm

If the spectral norm is used to prescribe the desired accuracy, we need to proceed differently. In this case, it is always possible to approximate very efficiently the spectral norm of the full matrix \mathbf{A} using few iterations, say k , of the power method (see e.g. [GL89]), using $\mathcal{O}(mnk)$ flops in total (if \mathbf{A} is stored as a full matrix).

However, it is difficult to estimate the absolute error

$$\|\mathbf{A} - \mathbf{W}\mathbf{Z}^\top\| = \|(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top)\mathbf{A}\|$$

without increasing the complexity of the algorithm. Indeed, such an approximation needs to be computed approximately r times (once for every new sample used to update \mathbf{W}), and using the power method each time would lead to the increased complexity $\mathcal{O}(mnrk)$.

We may tackle this issue by adding more than one random sample at a time, as described in Algorithm 4: given the sampling step $q > 0$, we only need to apply the power method r/q times, leading to an overall complexity of $\mathcal{O}(mnr(1 + k/q))$.

Algorithm 4: Randomized compression with fixed (spectral) precision with the power method

Data: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\varepsilon > 0$, integer $q > 0$

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ and such that (3.8) holds for the spectral norm

```

1 Compute  $\|\mathbf{A}\|$ 
2  $\mathbf{W} := []$ ,  $\mathbf{Z} := []$ 
3 while  $\|(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top)\mathbf{A}\| > \varepsilon\|\mathbf{A}\|$  do
4   for  $l = 1, \dots, q$  do
5     Generate Gaussian vector  $\boldsymbol{\omega} \in \mathbb{R}^n$ 
6      $\mathbf{w} := (\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top)\mathbf{A}\boldsymbol{\omega}$ 
7      $\mathbf{w} := \mathbf{w}/\|\mathbf{w}\|$ 
8      $\mathbf{W} := [\mathbf{W} \mid \mathbf{w}]$ 
9      $\mathbf{z} := \mathbf{A}^\top \mathbf{w}$ 
10     $\mathbf{Z} := [\mathbf{Z} \mid \mathbf{z}]$ 
11  end
12  Compute  $\|(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top)\mathbf{A}\|$ 
13 end
```

Error estimation with a randomized approach

A different approach, which does not rely on a staggered sampling of the error, can be derived as described in [HMT11]: instead of computing the actual value of $\|(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top)\mathbf{A}\|$, we approximate it using a randomized estimator, which is much cheaper to evaluate.

In particular, we can estimate the spectral norm of a matrix by using samples from its range, as shown by the following result, which is a generalization of Lemma 4.1 in [HMT11] and Theorem 1 in [Dix83].

Theorem 3.3. Let $\mathbf{B} \in \mathbb{R}^{m \times n}$, and consider $q > 0$ i.i.d. random vectors in $\{\omega^{(l)}\}_{l=1}^q \subset \mathbb{R}^n$, which follow a unitarily invariant distribution \mathcal{D} . Then, for any $\theta > 0$ it holds that

$$\|\mathbf{B}\| \leq \theta \max_{l=1, \dots, q} \|\mathbf{B}\omega^{(l)}\| \quad (3.13)$$

except with probability $\Pr(|\omega| < \frac{1}{\theta})^q$ at most, with ω being a random variable distributed as one of the components of $\omega^{(1)}$.

Proof. We follow the idea of the proof of Theorem 1 in [Dix83].

If $\|\mathbf{B}\| = 0$, the claim is trivial. Assume $\|\mathbf{B}\| > 0$ and let $l \in \{1, \dots, q\}$. We consider the (full) SVD of \mathbf{B} :

$$\mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$$

and compute the components of $\omega^{(l)}$ along the right singular vectors of \mathbf{B} :

$$\boldsymbol{\xi} = \mathbf{V}^\top \omega^{(l)} \quad \text{i.e.} \quad \omega^{(l)} = \mathbf{V}\boldsymbol{\xi}$$

Since \mathcal{D} is unitarily invariant, $\boldsymbol{\xi}$ still follows the same distribution as $\omega^{(l)}$.

We can express the norm of $\mathbf{B}\omega^{(l)}$ as

$$\|\mathbf{B}\omega^{(l)}\| = \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \omega^{(l)}\| = \|\mathbf{U}\boldsymbol{\Sigma}\boldsymbol{\xi}\| = \|\boldsymbol{\Sigma}\boldsymbol{\xi}\| = \left(\sum_{j=1}^n \xi(j)^2 \Sigma(j, j)^2 \right)^{1/2}$$

which implies that

$$\|\mathbf{B}\omega^{(l)}\| \geq |\xi(1)| \Sigma(1, 1) = |\xi(1)| \|\mathbf{B}\|$$

Hence, the event $\{\|\mathbf{B}\| > \theta \|\mathbf{B}\omega^{(l)}\|\}$ implies $\{\|\mathbf{B}\| > \theta |\xi(1)| \|\mathbf{B}\|\}$, and we can derive a bound on the failure probability

$$\begin{aligned} \Pr\left(\|\mathbf{B}\| > \theta \|\mathbf{B}\omega^{(l)}\|\right) &\leq \Pr\left(\|\mathbf{B}\| > \theta |\xi(1)| \|\mathbf{B}\|\right) \\ &= \Pr\left(1 > \theta |\xi(1)|\right) \\ &= \Pr\left(|\omega| < \frac{1}{\theta}\right) \end{aligned}$$

with ω being a random variable distributed as $\boldsymbol{\xi}(1)$.

Using the independence of the samples $\{\omega^{(l)}\}_{l=1}^q$, we derive

$$\begin{aligned} \Pr\left(\|\mathbf{B}\| > \theta \max_{l=1, \dots, q} \|\mathbf{B}\omega^{(l)}\|\right) &= \Pr\left(\bigcap_{l=1}^q \{\|\mathbf{B}\| > \theta \|\mathbf{B}\omega^{(l)}\|\}\right) = \\ &= \prod_{l=1}^q \Pr\left(\|\mathbf{B}\| > \theta \|\mathbf{B}\omega^{(l)}\|\right) \leq \Pr\left(|\omega| < \frac{1}{\theta}\right)^q \end{aligned}$$

which, after taking the complement, implies the claim. \square

In particular, since we have chosen to apply Gaussian vectors to solve the compression problem, we can restrict Theorem 3.3 to the Gaussian case (which is unitarily invariant) and obtain the following corollary.

Corollary 3.3.1. Let $\mathbf{B} \in \mathbb{R}^{m \times n}$, and consider $q > 0$ independent Gaussian random vectors $\{\boldsymbol{\omega}^{(l)}\}_{l=1}^q \subset \mathbb{R}^n$. Then, for any $\alpha > 0$ it holds that

$$\|\mathbf{B}\| \leq \sqrt{\frac{2}{\pi}} \alpha \max_{l=1, \dots, q} \|\mathbf{B}\boldsymbol{\omega}^{(l)}\| \quad (3.14)$$

except with probability α^{-q} at most.

Proof. If $\alpha \leq 1$, the claim holds trivially, since the bound on the failure probability is at least 1.

Let $\alpha > 1$. We define $\theta > 0$ as the unique value which satisfies

$$\alpha^{-1} = \Pr\left(|\omega| < \frac{1}{\theta}\right)$$

with ω being a 1-dimensional Gaussian random variable. By applying Theorem 3.3 with this choice of θ , we obtain

$$\Pr\left(\|\mathbf{B}\| \leq \theta \max_{l=1, \dots, q} \|\mathbf{B}\boldsymbol{\omega}^{(l)}\|\right) \geq 1 - \alpha^{-q}$$

Now, we observe that

$$\alpha^{-1} = \int_{-1/\theta}^{1/\theta} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \leq \int_{-1/\theta}^{1/\theta} \frac{1}{\sqrt{2\pi}} dx = \sqrt{\frac{2}{\pi}} \theta^{-1}$$

which implies that $\theta \leq \sqrt{\frac{2}{\pi}} \alpha$.

Thus, it holds

$$\Pr\left(\|\mathbf{B}\| \leq \sqrt{\frac{2}{\pi}} \alpha \max_{l=1, \dots, q} \|\mathbf{B}\boldsymbol{\omega}^{(l)}\|\right) \geq \Pr\left(\|\mathbf{B}\| \leq \theta \max_{l=1, \dots, q} \|\mathbf{B}\boldsymbol{\omega}^{(l)}\|\right) \geq 1 - \alpha^{-q} \quad \square$$

Using this result, we can devise a very efficient procedure to compress a matrix with fixed precision (in the spectral norm), on the wake of Algorithm 4.2 in [HMT11]. Assume that $\mathbf{W} \in \mathbb{R}^{m \times k}$ is the current approximation (with orthonormal columns) of the left factor, so that $\mathbf{A} \approx \mathbf{W}\mathbf{W}^\top \mathbf{A}$. To estimate the current error, we consider $q > 0$ independent Gaussian vectors $\{\boldsymbol{\omega}^{(l)}\}_{l=1}^q \subset \mathbb{R}^n$ and we compute

$$10\sqrt{\frac{2}{\pi}} \max_{l=1, \dots, q} \|(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A}\boldsymbol{\omega}^{(l)}\|$$

If the result is smaller than $\varepsilon \|\mathbf{A}\|$, we stop the algorithm, since the desired precision ε is achieved (with probability at least $1 - 10^{-q}$). Otherwise, we add a new column to \mathbf{W} and repeat the check.

We summarize this procedure in Algorithm 5.

To evaluate the complexity of the algorithm, we have to take into account the fact that the set of samples $\{(\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A}\boldsymbol{\omega}^{(l)}\}_{l=1}^q$ is updated recursively. Indeed, each step requires the computation of a single new sample from the range of \mathbf{A} , and its orthogonalization with respect to the range of \mathbf{W} (see line 14), as well as the update of the right factor (see line 12). Both operations take $\mathcal{O}(mn)$ flops, and are repeated r

Algorithm 5: Randomized compression with fixed (spectral) precision with a randomized norm estimate

Data: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\varepsilon > 0$, integer $q > 0$

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ and such that (3.8) holds with high probability for the spectral norm

```

1 Compute  $\|\mathbf{A}\|$ 
2 for  $l = 1, \dots, q$  do
3   | Generate Gaussian vector  $\boldsymbol{\omega} \in \mathbb{R}^n$ 
4   |  $\mathbf{y}^{(l)} := \mathbf{A}\boldsymbol{\omega}$ 
5 end
6  $\mathbf{W} := []$ ,  $\mathbf{Z} := []$ 
7 while  $10\sqrt{2/\pi} \max_{l=1, \dots, q} \|\mathbf{y}^{(l)}\| \geq \varepsilon \|\mathbf{A}\|$  do
8   |  $k :=$  random element of  $\{1, \dots, q\}$ 
9   |  $\mathbf{w} := \mathbf{y}^{(k)} / \|\mathbf{y}^{(k)}\|$ 
10  |  $\mathbf{W} := [\mathbf{W} \mid \mathbf{w}]$ 
11  |  $\mathbf{z} := \mathbf{A}^\top \mathbf{w}$ 
12  |  $\mathbf{Z} := [\mathbf{Z} \mid \mathbf{z}]$ 
13  | Generate Gaussian vector  $\boldsymbol{\omega} \in \mathbb{R}^n$ 
14  |  $\mathbf{y}^{(k)} := (\mathbf{I}_m - \mathbf{W}\mathbf{W}^\top) \mathbf{A}\boldsymbol{\omega}$ 
15  | for  $l = 1, \dots, k-1, k+1, \dots, q$  do
16  |   |  $\mathbf{y}^{(l)} := (\mathbf{I}_m - \mathbf{w}\mathbf{w}^\top) \mathbf{y}^{(l)}$ 
17  |   end
18 end
```

times. If we include the cost of the initialization of the set of samples (which requires $\mathcal{O}(mn)$ flops per sample), the overall complexity of the algorithm is $\mathcal{O}(mn(r+q))$.

Also, as in the case of the Frobenius norm, it may be desirable to add a recompression step (which costs $\mathcal{O}(m(r+q)^2 + nr(r+q))$ flops) at the end of the procedure.

3.3 Numerical examples

3.3.1 Computational environment

All numerical experiments described in this and in the upcoming chapters refer to implementations in MATLAB[®] 2016a, executed on the Malc1 cluster at EPFL (malc1.epfl.ch). In particular, we use 8 cores belonging to a single node with an Intel[®] Xeon[®] X5675E processor with frequency 3.07 GHz, and 191 GB of RAM.

Moreover, the MATLAB[®] implementations of the algorithms described in this thesis are available at https://github.com/pradover/TT_HOSVD.

3.3.2 Discretization of a rational function

Given the integer $n > 0$, we consider a uniform discretization of the interval $[\frac{1}{10}, \frac{n}{10}]$: $x_i = \frac{i}{10}$ for $i \in \{1, \dots, n\}$. For any fixed $\gamma \in \mathbb{R}^+$, we define the matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ entry-wise as

$$\mathbf{A}(i, j) = f_{2, \gamma}(x_i, x_j) = (x_i + x_j)^{-\gamma} \quad \text{for } (i, j) \in \{1, \dots, n\}^2$$

Due to the smoothness of $f_{2,\gamma}$ in $(0, \frac{n+1}{10})^2$, the singular values of \mathbf{A} decay with (at least) exponential rate (see e.g. [GL17]). For example, if $\gamma = 1$ and $n = 10^5$, the numerical rank is only 54 (i.e. there are exactly 54 singular values larger than the machine epsilon).

We consider the choices of γ and n specified above, and we apply Algorithm 1 to compress \mathbf{A} down to rank $r \in \{10, 20, 30\}$. In particular, we wish to investigate how the solution depends on the oversampling parameter $p \in \{0, 1, \dots, 8\}$.

In Figure 3.1 we report the relative error (in the Frobenius norm) of each low-rank approximation. In particular, the plots show the average, minimum and maximum values of the error after $N = 20$ simulation with each configuration (r, p) . It appears that small values of p (e.g. $p = 4$) are enough to achieve the optimal accuracy, in agreement with the numerical observations in [HMT11].

Moreover, we can verify that the theoretical bounds of Section 3.1.1 are satisfied. Indeed, the average error satisfies the bound (3.4) (see green dashed line), and (3.6) holds for all samples (see black dash-dotted line).

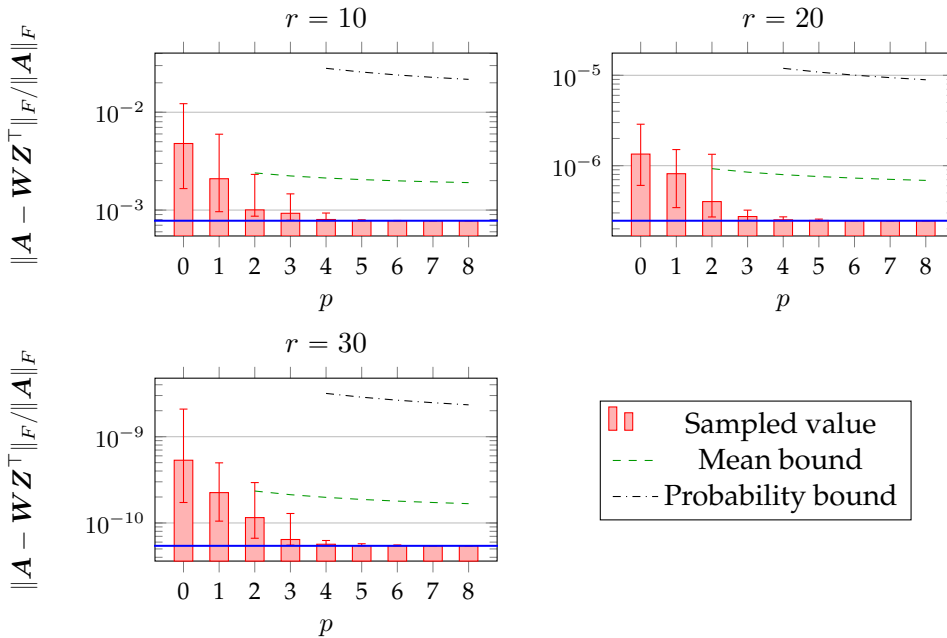


FIGURE 3.1: Average relative error versus oversampling for Algorithm 1. The error bars indicate the range of values obtained over $N = 20$ simulations. In blue the optimal error. In green and black the bounds from Theorems 3.1 and 3.2 respectively.

3.3.3 Compression of a Matérn kernel

Given the two positive parameters ν and ρ , we define the Matérn covariance function:

$$C_\nu^\rho(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\rho} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}r}{\rho} \right)$$

for $r \geq 0$, with K_ν being the modified Bessel function of the second kind. Such functions are employed in several machine learning techniques, among which we cite kernel regression (see e.g. [Bis09]) and Gaussian processes (see e.g. [RW05]).

Given the integer $n > 0$, we consider a uniform discretization of the interval $[\frac{1}{2n}, 1 - \frac{1}{2n}]$ ($x_i = \frac{2i-1}{2n}$, for $i \in \{1, \dots, n\}$), and we define the symmetric matrix

$\mathbf{A} \in \mathbb{R}^{n \times n}$ as

$$\mathbf{A}(i, j) = C_\nu^\rho (|x_i - x_j|) \quad \text{for } (i, j) \in \{1, \dots, n\}^2$$

We fix $\rho = 1$, $\nu = \frac{3}{2}$ and $n = 10^5$, so that the singular values of \mathbf{A} have a slow polynomial decay ($\sigma_j \sim j^{-4.3}$ approximately). Our objective is to compress the matrix with fixed precision $\varepsilon \in \{10^{-2}, 10^{-6}\}$ in the Frobenius norm. To this aim we apply Algorithm 3 (including the recompression step) with $q \in \{1, 2, 3, 5, 7, 10, 15, 20\}$, and observe the results of $N = 20$ trials with each configuration (ε, q) .

The relative error of the approximation is shown in Figure 3.2. For $\varepsilon = 10^{-2}$ we can see that really small values of q (even $q = 1$) are sufficient to achieve the desired precision. However, for $\varepsilon = 10^{-6}$ the average error is smaller than ε only for $q \geq 15$.

This is due to the error introduced by the approximation (3.12). Indeed, since the singular values of \mathbf{A} decay quite slowly, the norm of $\tilde{\mathbf{Z}}$ does not approximate the norm of the full matrix with enough accuracy for small values of q . Thus, Algorithm 3 is stopped before a sufficiently rich approximation of the range of \mathbf{A} is identified.

In Figure 3.3 we report the rank of the approximation obtained, both before and after the recompression step. For both choices of ε , we can observe that the discrepancy between the ranks before and after recompression is approximately equal to q . This is due to the fact that (3.12) is consistent with the condition used in the recompression step, since both rely on $\tilde{\mathbf{Z}}$ to estimate the norm of the full matrix.

For $\varepsilon = 10^{-2}$, the final ranks always equal the optimal one, attesting that the approximate condition (3.12) is reliable in stopping the algorithm. This is not true for $\varepsilon = 10^{-6}$, where the final ranks are below the optimal one for small values of q .

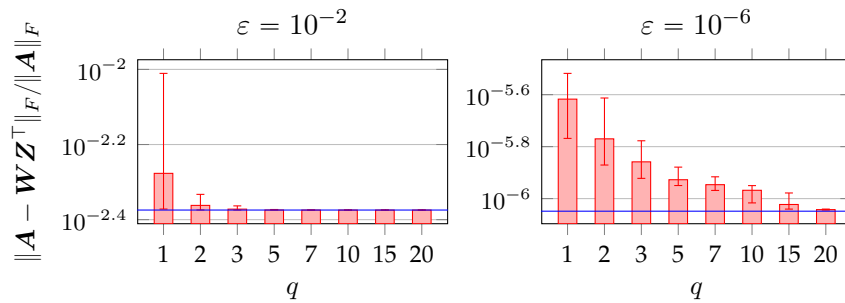


FIGURE 3.2: Average relative error versus sampling parameter for Algorithm 3. The error bars indicate the range of values obtained over $N = 20$ simulations. In blue the optimal error.

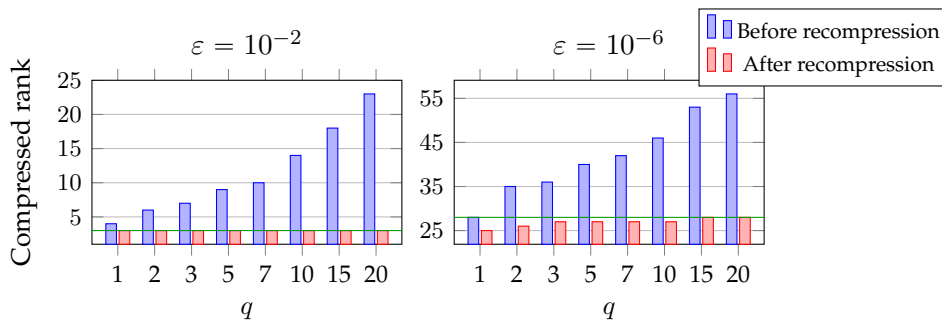


FIGURE 3.3: Modal values of the compressed rank (before and after recompression) versus sampling parameter for Algorithm 3. In green the optimal rank.

3.3.4 Compression with non-Gaussian random vectors

We consider the matrix A defined in Section 3.3.3. Our aim is to compress it to ranks $r \in \{5, 15, 30\}$ using modified versions of Algorithm 1, where the random matrix Ω is sampled according to different distributions. In particular, we consider sample vectors with independent entries drawn from the uniform distribution over $(0, 1)$, and “rank-1” Gaussian sample vectors, i.e. vectors of the form $\omega = \omega_1 \otimes \omega_2$, with ω_1 and ω_2 being \sqrt{n} -dimensional (independent) Gaussian vectors.

For these distributions there are no theoretical results similar to the ones presented in Section 3.1.1. However, they have advantages from the computational point of view: both are (on most calculators) cheaper to sample than the Gaussian distribution. Moreover, in the “rank-1” case, the structured format of the vector can be exploited to gain efficiency (see e.g. [ZJD15], [KP16], and Sections 4.3 and 5.2 of this thesis).

In Figure 3.4 we compare the relative errors with the ones achieved in the Gaussian case. For all approaches, we choose $p \in \{1, 2, 3, 5, 7, 10, 15, 20\}$ as the sampling parameter.

We can observe that the results obtained with the three distributions are similar: all of them appear to converge at the same rate, and the ranges of the values are of similar size. In particular, all appear to satisfy (3.4) and (3.6).

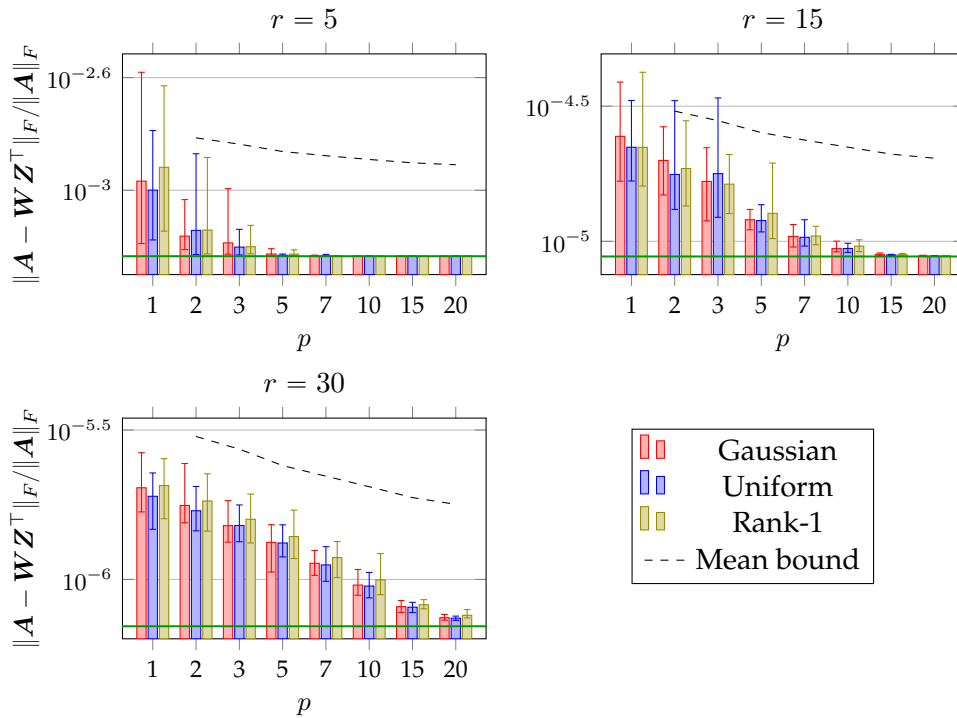


FIGURE 3.4: Average relative error versus oversampling for Algorithm 1. The bars connect the minimum and the maximum value of the relative error obtained over $N = 20$ simulations. In green the optimal error. In black the bound from Theorem 3.1. The bound from Theorem 3.2 is not shown for ease of visualization, since it is several orders of magnitude larger than the sampled values.

Chapter 4

Compression of tensors in the TT format

The Tensor Train (TT) format was introduced by Oseledets and Tyrtyshnikov (see [OT09]) as a way to break the curse of dimensionality. Indeed, most operations in the TT format have polynomial complexity in the order of the tensor. The format had already been used for a few decades in several physics applications (see e.g. [Aff+87] and [FNW92]) under the name Matrix Product States (MPS).

4.1 Overview of the Tensor Train format

We say that the d -dimensional tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ admits a TT representation with TT-ranks $(r_0 = 1, r_1, \dots, r_{d-1}, r_d = 1)$ if there exist d tensors of order 3

$$\mathcal{A}_j \in \mathbb{R}^{r_{j-1} \times n_j \times r_j} \quad \text{for } j \in \{1, \dots, d\}$$

(which are commonly called *cores* of the representation), such that

$$\mathcal{A}(i_1, \dots, i_d) = \sum_{\alpha_1, \dots, \alpha_{d-1}} \mathcal{A}_1(1, i_1, \alpha_1) \mathcal{A}_2(\alpha_1, i_2, \alpha_2) \dots \mathcal{A}_d(\alpha_{d-1}, i_d, 1) \quad (4.1)$$

for $(i_1, \dots, i_d) \in I_1 \times \dots \times I_d$. In the expression above, the sum over α_j is carried out from 1 to r_j (for $j \in \{1, \dots, d-1\}$), and we define

$$I_j = \{1, \dots, r_j\} \quad \text{for } j \in \{1, \dots, d\}$$

We will use the following notation to represent the mode-2 slices of each core

$$\mathcal{A}_j(i_j) = \mathcal{A}_j(:, i_j, :) \in \mathbb{R}^{r_{j-1} \times r_j} \quad \text{for } i_j \in I_j, \text{ for } j \in \{1, \dots, d\},$$

so that we can interpret (4.1) as a chain of matrix products:

$$\mathcal{A}(i_1, \dots, i_d) = \mathcal{A}_1(i_1) \mathcal{A}_2(i_2) \dots \mathcal{A}_d(i_d) \quad \text{for } (i_1, \dots, i_d) \in I_1 \times \dots \times I_d \quad (4.2)$$

Another useful (graphic) representation of tensors in the TT format is given by *tensor network diagrams* (see e.g. [Oru14] and [Ste16]). An example of such a representation is given in Figure 4.1: each node corresponds to a core tensor, and each edge to an index over which a sum is computed.

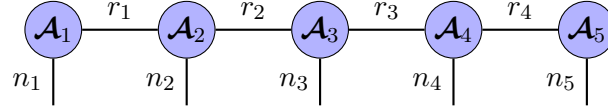


FIGURE 4.1: Tensor network diagram for a 5-dimensional tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_5}$ in the TT format (4.2) with TT-ranks $(1, r_1, r_2, r_3, r_4, 1)$.

For fixed TT-ranks, the TT representation of a tensor is not unique, similarly to the non-uniqueness of the dyadic representation with fixed rank (1.1) of a matrix. However, the following result guarantees the existence of a TT decomposition and an upper-bound on the minimal TT-ranks. Moreover, it provides an algorithm to convert any tensor into the TT format.

Theorem 4.1. *Let $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$. There exists a decomposition of the form (4.2) with $r_0 = r_d = 1$ and*

$$r_j = \text{rank } \mathcal{A}^{<j>} \quad \text{for } j \in \{1, \dots, d-1\} \quad (4.3)$$

Proof. We follow the idea of the proof of Theorem 2.1 in [Ose11].

Define (r_1, \dots, r_{d-1}) as in (4.3). Consider a rank- r_1 dyadic decomposition of $\mathcal{A}^{<1>}$:

$$\mathcal{A}^{<1>} = \mathbf{W}_1 \mathbf{Z}_1^\top \quad (4.4)$$

with

$$\mathbf{W}_1 \in \mathbb{R}^{n_1 \times r_1} \quad \text{and} \quad \mathbf{Z}_1 \in \mathbb{R}^{(n_2 \dots n_d) \times r_1}$$

We define the first core $\mathcal{A}_1 \in \mathbb{R}^{1 \times n_1 \times r_1}$ so that $(\mathcal{A}_1)^{<2>} = \mathbf{W}_1$.

Let \mathcal{B}_1 be the d -dimensional tensor of size $r_1 \times n_2 \times n_3 \times \dots \times n_d$ such that

$$(\mathcal{B}_1)^{<1>} = \mathbf{Z}_1^\top \quad (4.5)$$

and assume $\tilde{r}_2 = \text{rank } (\mathcal{B}_1)^{<2>}$.

Consider a rank- \tilde{r}_2 dyadic decomposition of $(\mathcal{B}_1)^{<2>}$:

$$(\mathcal{B}_1)^{<2>} = \mathbf{W}_2 \mathbf{Z}_2^\top \quad (4.6)$$

with

$$\mathbf{W}_2 \in \mathbb{R}^{(r_1 n_2) \times \tilde{r}_2} \quad \text{and} \quad \mathbf{Z}_2 \in \mathbb{R}^{(n_3 \dots n_d) \times \tilde{r}_2}$$

We define the second core $\mathcal{A}_2 \in \mathbb{R}^{r_1 \times n_2 \times \tilde{r}_1}$ so that $(\mathcal{A}_2)^{<2>} = \mathbf{W}_2$. Then we proceed recursively to compute all the cores of the TT representation.

At the end of this left-to-right sweep, the TT-ranks of the decomposition are $(1, r_1, \tilde{r}_2, \dots, \tilde{r}_{d-1}, 1)$, where $\tilde{r}_j = \text{rank } (\mathcal{B}_{j-1})^{<2>}$ for $j \in \{2, \dots, d-1\}$. Hence, it suffices to prove that

$$\text{rank } (\mathcal{B}_{j-1})^{<2>} = \tilde{r}_j \leq r_j = \text{rank } \mathcal{A}^{<j>} \quad \text{for } j \in \{2, \dots, d-1\} \quad (4.7)$$

To this aim, we prove that

$$\text{rank } (\mathcal{B}_1)^{<l>} \leq \text{rank } \mathcal{A}^{<l>} \quad \text{for } l \in \{2, \dots, d-1\} \quad (4.8)$$

Proof of claim (4.8): From (4.4) we deduce that

$$\mathbf{Z}_1^\top = (\mathbf{W}_1^\top \mathbf{W}_1)^{-1} \mathbf{W}_1^\top \mathcal{A}^{<1>} =: \widetilde{\mathbf{W}}_1^\top \mathcal{A}^{<1>}$$

or equivalently

$$\mathbf{Z}_1^\top = \left(\mathcal{A} \times_1 \widetilde{\mathbf{W}}_1^\top \right)^{<1>}$$

Then (4.5) implies that $\mathbf{B}_1 = \mathcal{A} \times_1 \widetilde{\mathbf{W}}_1^\top$.

Now, let $l \in \{2, \dots, d-1\}$ be fixed and consider the rank- r_l dyadic decomposition of $\mathcal{A}^{<l>}$:

$$\mathcal{A}^{<l>} = \mathbf{W}'_l \mathbf{Z}'_l{}^\top$$

It is easy to verify that the following rank- r_l decomposition holds for $(\mathbf{B}_1)^{<l>}$:

$$(\mathbf{B}_1)^{<l>} = \left(\mathbf{I}_{n_l} \otimes \dots \otimes \mathbf{I}_{n_2} \otimes \widetilde{\mathbf{W}}_1^\top \right) \mathbf{W}'_l \mathbf{Z}'_l{}^\top = \mathbf{W}''_l \mathbf{Z}'_l{}^\top$$

This implies that $\text{rank} (\mathbf{B}_1)^{<l>} \leq r_l$ and proves the claim. \square

Using a similar argument (by replacing \mathcal{A} with \mathbf{B}_{j-1} and \mathbf{B}_1 with \mathbf{B}_j), we can verify that

$$\text{rank} (\mathbf{B}_j)^{<l>} \leq \text{rank} (\mathbf{B}_{j-1})^{<l+1>} \quad \text{for } l \in \{2, \dots, d-j\}$$

for $j \in \{2, \dots, d-2\}$. The claim (4.7) follows immediately, since

$$\text{rank} (\mathbf{B}_{j-1})^{<2>} \leq \text{rank} (\mathbf{B}_{j-2})^{<3>} \leq \dots \leq \text{rank} (\mathbf{B}_1)^{<j>} \leq \text{rank} \mathcal{A}^{<j>}$$

for $j \in \{2, \dots, d-1\}$.

At this point, we have obtained a TT decomposition whose j^{th} TT-rank is at most r_j (for $j \in \{1, \dots, d-1\}$). This proves the theorem, since, in order to obtain a TT decomposition with TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$, it suffices to add a zero-padding to all the cores with deficient size. \square

The proof of the previous Theorem is constructive. As such, it provides an algorithm to convert any tensor into the TT format. The complexity of the algorithm is dominated by the cost of the computation of the d dyadic decompositions of the form (4.4) and (4.6), which can be carried out using a RRQR decomposition: if we assume for simplicity that

$$n_1 = \dots = n_d = n \quad \text{and} \quad r_1 = \dots = r_{d-1} = r,$$

the cost for computing $\mathcal{A}^{<1>}$ is $\mathcal{O}(n^d r)$ flops, while for $(\mathbf{B}_j)^{<2>}$ (with $j \in \{1, \dots, d-2\}$) it is $\mathcal{O}(n^{d-j} r^2)$ flops, leading to an overall complexity of $\mathcal{O}(n^d r^2)$.

The complexity of the conversion to the TT format with respect to d can be reduced to polynomial time if the initial tensor \mathcal{A} is stored in some structured format. For instance, the cost is linear in d if \mathcal{A} is stored using the canonical decomposition (see

e.g. [Ose11]):

$$\mathcal{A}(i_1, \dots, i_d) = \sum_{\alpha=1}^{\bar{r}} U_1(i_1, \alpha) \dots U_d(i_d, \alpha) \quad \text{for } (i_1, \dots, i_d) \in I_1 \times \dots \times I_d \quad (4.9)$$

with \bar{r} being the canonical rank of \mathcal{A} , and $U_j \in \mathbb{R}^{n_j \times \bar{r}}$ (for $j \in \{1, \dots, d\}$) being the canonical factors.

It is easy to compute a TT-decomposition of \mathcal{A} with TT-ranks $(1, \bar{r}, \dots, \bar{r}, 1)$ by introducing $d - 2$ additional identity matrices of size $\bar{r} \times \bar{r}$ into (4.9): for any $(i_1, \dots, i_d) \in I_1 \times \dots \times I_d$, it holds

$$\mathcal{A}(i_1, \dots, i_d) = \sum_{\alpha_1, \dots, \alpha_{d-1}} U_1(i_1, \alpha_1) \mathbf{I}_{\bar{r}}(\alpha_1, \alpha_2) \dots \mathbf{I}_{\bar{r}}(\alpha_{d-2}, \alpha_{d-1}) U_{d-1}(i_{d-1}, \alpha_{d-1}) U_d(i_d, \alpha_{d-1})$$

where the sum ranges from 1 to \bar{r} for each α_j (for $j \in \{1, \dots, d-1\}$).

Now we can define the inner cores of the TT representation:

$$\mathcal{A}_j(\alpha_{j-1}, i_j, \alpha_j) = \mathbf{I}_{\bar{r}}(\alpha_{j-1}, \alpha_j) U_j(i_j, \alpha_j) \quad \text{for } (i_j, \alpha_{j-1}, \alpha_j) \in I_j \times \{1, \dots, \bar{r}\}^2$$

for $j \in \{2, \dots, d-1\}$, and the outer cores accordingly:

$$\mathcal{A}_1(1, i_1, \alpha_1) = U_1(i_1, \alpha_1) \quad \text{for } (i_1, \alpha_1) \in I_1 \times \{1, \dots, \bar{r}\}$$

$$\mathcal{A}_d(\alpha_{d-1}, i_d, 1) = U_d(i_d, \alpha_{d-1}) \quad \text{for } (i_d, \alpha_{d-1}) \in I_d \times \{1, \dots, \bar{r}\}$$

4.2 Compression of tensors in the TT format using deterministic methods

Now we focus on the following problem: let \mathcal{A} be a tensor in the TT format (4.2) with TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$. We want to find a new (possibly approximate) TT decomposition of \mathcal{A} with cores \mathcal{A}_j^* whose TT-ranks $(1, r_1^*, \dots, r_{d-1}^*, 1)$ are smaller.

Two interesting formulations of this problem can be stated in same spirit as in Chapter 3:

- compress the tensor down to some arbitrary TT-ranks;
- compress the tensor so that the relative Frobenius error between the original tensor and the new one \mathcal{A}^* , i.e.

$$\frac{\|\mathcal{A}^* - \mathcal{A}\|_F}{\|\mathcal{A}\|_F}, \quad (4.10)$$

is smaller than some prescribed accuracy.

Using the terminology from [Ose11], we will refer to the two cases as *truncation* and *rounding* problems respectively.

To devise an efficient algorithm to approach these problems, it is useful to define the concept of left- and right-orthogonality for a tensor in the TT format.

4.2.1 Orthogonalization of tensors in the TT format

Being the natural generalization of the QR decomposition for tensors in the TT format, orthogonalization inherits a crucial role in the manipulation of tensors.

For ease of notation, we first introduce the concept of interface matrices for tensors in the TT format.

Definition 4.1 (Interface matrices). *Let \mathcal{A} be a d -dimensional tensor in the TT format (4.2) with TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$ and let $j \in \{1, \dots, d-1\}$. We consider the $(j+1)$ -dimensional tensor \mathcal{B} of size $n_1 \times \dots \times n_j \times r_j$ with entries*

$$\mathcal{B}(i_1, \dots, i_j, :) = \mathbf{A}_1(i_1) \dots \mathbf{A}_j(i_j) \quad \text{for } (i_1, \dots, i_j) \in I_1 \times \dots \times I_j$$

and we define the left interface matrix as

$$\mathcal{A}_{\leq j} = \mathcal{B}^{<j>} \in \mathbb{R}^{(n_1 \dots n_j) \times r_j} \quad (4.11)$$

Similarly, we consider the $(d-j+1)$ -dimensional tensor \mathcal{C} of size $r_j \times n_{j+1} \times \dots \times n_d$ with entries

$$\mathcal{C}(:, i_{j+1}, \dots, i_d) = \mathbf{A}_{j+1}(i_{j+1}) \dots \mathbf{A}_d(i_d) \quad \text{for } (i_{j+1}, \dots, i_d) \in I_{j+1} \times \dots \times I_d$$

and we define the right interface matrix as

$$\mathcal{A}_{\geq j+1} = (\mathcal{C}^{<1>})^\top \in \mathbb{R}^{(n_{j+1} \dots n_d) \times r_j} \quad (4.12)$$

Since the TT factorization of a tensor is not unique, it is important to observe that interface matrices depend on the particular TT decomposition and not just on the tensor itself.

The importance of interface matrices is due to the fact that they appear in the dyadic decompositions of tensor unfoldings: given $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format (4.2), for $j \in \{1, \dots, d-1\}$, it holds

$$\mathcal{A}^{<j>} = \mathcal{A}_{\leq j} \mathcal{A}_{\geq j+1}^\top$$

Given this result, we define the orthogonality of a tensor in the TT format in terms of the orthogonality of its interface matrices.

Definition 4.2 (Left- and right-orthogonality). *Let \mathcal{A} be a d -dimensional tensor in the TT format (4.2) with TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$. We say that \mathcal{A} is left-orthogonal if*

$$\mathcal{A}_{\leq j}^\top \mathcal{A}_{\leq j} = \mathbf{I}_{r_j} \quad (4.13)$$

for $j \in \{1, \dots, d-1\}$, and that it is right-orthogonal if

$$\mathcal{A}_{\geq j+1}^\top \mathcal{A}_{\geq j+1} = \mathbf{I}_{r_j} \quad (4.14)$$

for $j \in \{1, \dots, d-1\}$.

There is an alternative way to characterize the orthogonality of a tensor, as stated by the following Lemma.

Lemma 4.2. Let \mathcal{A} be as in Definition 4.2. \mathcal{A} is left-orthogonal if and only if

$$(\mathcal{A}_j)^{\langle 2 \rangle \top} (\mathcal{A}_j)^{\langle 2 \rangle} = \mathbf{I}_{r_j} \quad (4.15)$$

for $j \in \{1, \dots, d-1\}$, and right-orthogonal if and only if

$$(\mathcal{A}_{j+1})^{\langle 1 \rangle} (\mathcal{A}_{j+1})^{\langle 1 \rangle \top} = \mathbf{I}_{r_j} \quad (4.16)$$

for $j \in \{1, \dots, d-1\}$.

Proof. We only provide the proof for the left-orthogonality. The other case can be treated in an analogous fashion.

For the first core we have

$$\mathcal{A}_{\leq 1} = (\mathcal{A}_1)^{\langle 2 \rangle}$$

Thus, (4.13) and (4.15) are equivalent for $j = 1$.

We proceed by induction. Let $l \in \{1, \dots, d-2\}$, and assume that (4.13) and (4.15) are equivalent for $j \in \{1, \dots, l\}$. We have that

$$\mathcal{A}_{\leq l+1} = (\mathbf{I}_{n_{l+1}} \otimes \mathcal{A}_{\leq l}) (\mathcal{A}_{l+1})^{\langle 2 \rangle}$$

and we can apply Proposition 2.2 to obtain

$$\begin{aligned} \mathcal{A}_{\leq l+1}^\top \mathcal{A}_{\leq l+1} &= (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathbf{I}_{n_{l+1}} \otimes \mathcal{A}_{\leq l})^\top (\mathbf{I}_{n_{l+1}} \otimes \mathcal{A}_{\leq l}) (\mathcal{A}_{l+1})^{\langle 2 \rangle} \\ &= (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathbf{I}_{n_{l+1}} \otimes (\mathcal{A}_{\leq l}^\top \mathcal{A}_{\leq l})) (\mathcal{A}_{l+1})^{\langle 2 \rangle} \end{aligned}$$

Assume that either (4.13) or (4.15) is satisfied for all values of $j \in \{1, \dots, d-1\}$. Then, in particular, we have that $\mathcal{A}_{\leq l}^\top \mathcal{A}_{\leq l} = \mathbf{I}_{r_l}$ due to the induction hypothesis. As such, it holds that

$$\begin{aligned} (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathbf{I}_{n_{l+1}} \otimes (\mathcal{A}_{\leq l}^\top \mathcal{A}_{\leq l})) (\mathcal{A}_{l+1})^{\langle 2 \rangle} &= (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathbf{I}_{n_{l+1}} \otimes \mathbf{I}_{r_l}) (\mathcal{A}_{l+1})^{\langle 2 \rangle} \\ &= (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} \mathbf{I}_{n_{l+1} r_l} (\mathcal{A}_{l+1})^{\langle 2 \rangle} \\ &= (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathcal{A}_{l+1})^{\langle 2 \rangle} \end{aligned}$$

which implies

$$\mathcal{A}_{\leq l+1}^\top \mathcal{A}_{\leq l+1} = (\mathcal{A}_{l+1})^{\langle 2 \rangle \top} (\mathcal{A}_{l+1})^{\langle 2 \rangle} \quad \square$$

Using this easier definition, we can efficiently orthogonalize a tensor \mathcal{A} in the TT format (either to the left or to the right) by following the intuition used in the proof of Theorem 4.1 (see e.g. [Ose11]).

We proceed as follows for the left-orthogonalization. We consider the QR decomposition of the second unfolding of the first core

$$(\mathcal{A}_1)^{\langle 2 \rangle} = \mathbf{Q}_1 \mathbf{R}_1$$

We overwrite the first core so that $(\mathcal{A}_1^*)^{\langle 2 \rangle} = \mathbf{Q}_1$, and we use \mathbf{R}_1 to update the second core:

$$\mathcal{A}_2^* = \mathcal{A}_2 \times_1 \mathbf{R}_1$$

In this way, (4.15) is satisfied for $j = 1$, and the overall tensor has not been affected, since

$$\mathbf{A}_1(i_1)\mathbf{A}_2(i_2) = \mathbf{Q}_1(i_1, :)\mathbf{R}_1\mathbf{A}_2(i_2) = \mathbf{A}_1^*(i_1)\mathbf{A}_2^*(i_2) \quad \text{for } (i_1, i_2) \in I_1 \times I_2$$

Now we consider $(\mathcal{A}_2^*)^{<2>}$ and proceed similarly, by computing its QR decomposition $(\mathcal{A}_2^*)^{<2>} = \mathbf{Q}_2\mathbf{R}_2$, overwriting the second core with \mathbf{Q}_2 (so that (4.15) is satisfied for $j = 2$), and updating the third core using \mathbf{R}_2 . Then we iterate over all the cores.

The resulting procedure is summarized in Algorithm 6. Assuming that

$$n_1 = \dots = n_d = n \quad \text{and} \quad r_1 = \dots = r_{d-1} = r,$$

the complexity of the algorithm is $\mathcal{O}(dnr^3)$ flops, since each QR decomposition, as well as each 1st mode (or 3rd mode) product, costs $\mathcal{O}(nr^3)$ flops.

We also provide a graphical depiction of the algorithm in Figure 4.2. In particular, we show the left-orthogonalization of a $n_1 \times \dots \times n_d$ tensor with TT-ranks $(1, r_1, r_2, r_3, 1)$. Each pink rectangle contains the two factors appearing in the QR decomposition of the gray block in the line above, after a suitable reshaping. It is interesting to observe that each of the matrices \mathbf{R}_j (with $j \in \{1, \dots, d-1\}$) represents an interface between a left-orthogonal and a non-orthogonal part, comprised respectively of green and blue nodes.

4.2.2 Truncation of tensors in the TT format in two sweeps

To solve the truncation problem, we can exploit the orthogonalization algorithm as follows. We right-orthogonalize the original tensor \mathcal{A} , leading to the TT decomposition with cores $\mathcal{A}'_1, \dots, \mathcal{A}'_d$.

We can express the first unfolding of \mathcal{A} using the newly obtained cores

$$\mathcal{A}^{<1>} = \mathcal{A}'_{\leq 1} \mathcal{A}'_{\geq 2}^\top = (\mathcal{A}'_1)^{<2>} \mathcal{A}'_{\geq 2}^\top$$

and compute the SVD of the left interface matrix

$$(\mathcal{A}'_1)^{<2>} = U_1 \Sigma_1 V_1^\top$$

Since \mathcal{A} has been right-orthogonalized, the right interface matrix $\mathcal{A}'_{\geq 2}$ has orthonormal columns. Hence the SVD of the first unfolding is of the form

$$\mathcal{A}^{<1>} = U_1 \Sigma_1 (\mathcal{A}'_{\geq 2} V_1)^\top$$

and it is possible to obtain the best rank- r_1^* approximation of $\mathcal{A}^{<1>}$ by truncating each of the factors of the SVD (see Chapter 3)

$$\mathbf{W}_1 = U_1(:, 1 : r_1^*) \quad \text{and} \quad \mathbf{Z}_1 = V_1(:, 1 : r_1^*) \Sigma_1(1 : r_1^*, 1 : r_1^*),$$

so that

$$(\mathcal{A}'_1)^{<2>} \approx \mathbf{W}_1 \mathbf{Z}_1^\top \quad \text{and} \quad \mathcal{A}^{<1>} \approx \mathbf{W}_1 (\mathcal{A}'_{\geq 2} \mathbf{Z}_1)^\top$$

Algorithm 6: Left- or right-orthogonalization of a tensor in the TT format**Data:** tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, direction**Result:** $\mathcal{A}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1^*, \dots, \mathcal{A}_d^*$, orthogonal in the specified direction

```

1 if direction=='left' then
2    $\mathcal{A}_1^* := \mathcal{A}_1$ 
3   for  $j = 1, 2, \dots, d-1$  do
4     Compute QR decomposition:  $(\mathcal{A}_j^*)^{<2>} =: Q_j R_j$ 
5      $\mathcal{A}_j^* := \text{reshape}(Q_j, [r_{j-1}, n_j, r_j])$ 
6      $\mathcal{A}_{j+1}^* := \mathcal{A}_{j+1} \times_1 R_j$ 
7   end
8 else
9    $\mathcal{A}_d^* := \mathcal{A}_d$ 
10  for  $j = d, d-1, \dots, 2$  do
11    Compute QR decomposition:  $(\mathcal{A}_j^*)^{<1>^\top} =: Q_j R_j$ 
12     $\mathcal{A}_j^* := \text{reshape}(Q_j^\top, [r_{j-1}, n_j, r_j])$ 
13     $\mathcal{A}_{j-1}^* := \mathcal{A}_{j-1} \times_3 R_j$ 
14  end
15 end

```

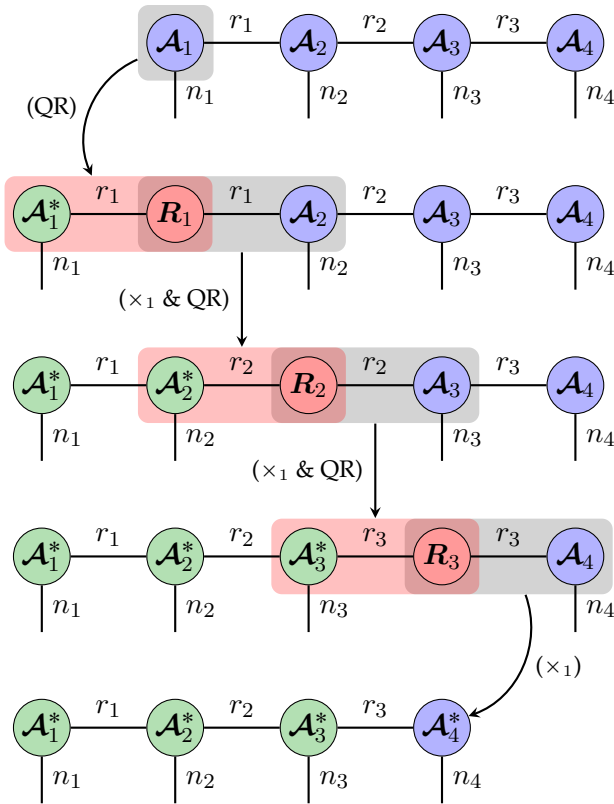


FIGURE 4.2: Tensor network diagram for the left-orthogonalization of $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_4}$, according to Algorithm 6. Green nodes have second unfoldings with orthonormal columns.

Just replacing the first core with its compressed version would not reduce the TT-ranks. Instead, we define the first core of \mathcal{A}^* so that $(\mathcal{A}_1^*)^{<2>} = \mathbf{W}_1$, and we compute a partial update of the second core:

$$\mathcal{A}_2^{**} = \mathcal{A}'_2 \times_1 \mathbf{Z}_1^\top$$

After these operations we have

$$\mathbf{A}'_1(i_1)\mathbf{A}'_2(i_2) \approx \mathbf{W}_1(i_1, :) \mathbf{Z}_1^\top \mathbf{A}'_2(i_2) = \mathcal{A}_1^*(i_1)\mathcal{A}_2^{**}(i_2) \quad \text{for } (i_1, i_2) \in I_1 \times I_2$$

with $\mathcal{A}_1^* \in \mathbb{R}^{1 \times n_1 \times r_1^*}$ and $\mathcal{A}_2^{**} \in \mathbb{R}^{r_1^* \times n_2 \times r_2}$.

Now we consider the second unfolding of \mathcal{A}

$$\begin{aligned} \mathcal{A}^{<2>} &= \mathcal{A}'_{\leq 2} \mathcal{A}'_{\geq 3}^\top = (\mathbf{I}_{n_2} \otimes \mathcal{A}'_{\leq 1}) (\mathcal{A}'_2)^{<2>} \mathcal{A}'_{\geq 3}^\top \\ &\approx (\mathbf{I}_{n_2} \otimes (\mathbf{W}_1 \mathbf{Z}_1^\top)) (\mathcal{A}'_2)^{<2>} \mathcal{A}'_{\geq 3}^\top \\ &= (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) (\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top) (\mathcal{A}'_2)^{<2>} \mathcal{A}'_{\geq 3}^\top \\ &= (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) (\mathcal{A}'_2 \times_1 \mathbf{Z}_1^\top)^{<2>} \mathcal{A}'_{\geq 3}^\top \\ &= (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) (\mathcal{A}_2^{**})^{<2>} \mathcal{A}'_{\geq 3}^\top \end{aligned}$$

and we compute the SVD of $(\mathcal{A}_2^{**})^{<2>}$: $(\mathcal{A}_2^{**})^{<2>} = \mathbf{U}_2 \mathbf{\Sigma}_2 \mathbf{V}_2^\top$.

We observe that

$$(\mathbf{I}_{n_2} \otimes \mathbf{W}_1)^\top (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) = \mathbf{I}_{n_2} \otimes (\mathbf{W}_1^\top \mathbf{W}_1) = \mathbf{I}_{n_2} \otimes \mathbf{I}_{r_1^*} = \mathbf{I}_{r_1^* n_2},$$

and that the right interface matrix $\mathcal{A}'_{\geq 3}$ has orthonormal columns. Hence, the SVD of the second unfolding is of the form

$$\mathcal{A}^{<2>} = ((\mathbf{I}_{n_2} \otimes \mathbf{W}_1) \mathbf{U}_2) \mathbf{\Sigma}_2 (\mathcal{A}'_{\geq 3} \mathbf{V}_2)^\top$$

Then, if we define

$$\mathbf{W}_2 = \mathbf{U}_2(:, 1 : r_2^*) \quad \text{and} \quad \mathbf{Z}_2 = \mathbf{V}_2(:, 1 : r_2^*) \mathbf{\Sigma}_2(1 : r_2^*, 1 : r_2^*),$$

we can compute the second core of \mathcal{A}^* so that $(\mathcal{A}_2^*)^{<2>} = \mathbf{W}_2$, and update the third core:

$$\mathcal{A}_3^{**} = \mathcal{A}'_3 \times_1 \mathbf{Z}_2^\top$$

Now we iterate over all the cores, leading to the procedure summarized in Algorithm 7. Assuming for simplicity that

$$n_1 = \dots = n_d = n, \quad r_1 = \dots = r_{d-1} = r \quad \text{and} \quad r_1^* = \dots = r_{d-1}^* = r^*,$$

the cost of each SVD is $\mathcal{O}(nr^3)$, while the cost of each 1st (or 3rd) mode product is $\mathcal{O}(nr^2r^*)$ flops. The algorithm could be sped up by replacing the SVD with a RRQR decomposition, which typically decreases the complexity to $\mathcal{O}(nr^2r^*)$. Still, the total cost cannot be smaller than the cost of the initial orthogonalization, i.e. $\mathcal{O}(dnr^3)$ flops.

Algorithm 7: Truncation of a tensor in the TT format

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, desired TT-ranks $(r_0^* = 1, r_1^*, \dots, r_{d-1}^*, r_d^* = 1)$

Result: left-orthogonal $\mathcal{A}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1^*, \dots, \mathcal{A}_d^*$ with the specified TT-ranks

```

1  $\mathcal{A}' := \text{Algorithm 6}(\mathcal{A}, \text{'right'})$ 
2  $\mathcal{A}_1^{**} := \mathcal{A}'_1$ 
3 for  $j = 1, 2, \dots, d - 1$  do
4   Compute SVD:  $(\mathcal{A}_j^{**})^{<2>} =: U_j \Sigma_j V_j^\top$ 
5    $W_j := U_j(:, 1:r_j^*)$ 
6    $Z_j := V_j(:, 1:r_j^*) \Sigma_j(1:r_j^*, 1:r_j^*)$ 
7    $\mathcal{A}_j^* := \text{reshape}(W_j, [r_{j-1}^*, n_j, r_j^*])$ 
8    $\mathcal{A}_{j+1}^{**} := \mathcal{A}'_{j+1} \times_1 Z_j^\top$ 
9 end
10  $\mathcal{A}_d^* := \mathcal{A}_d^{**}$ 

```

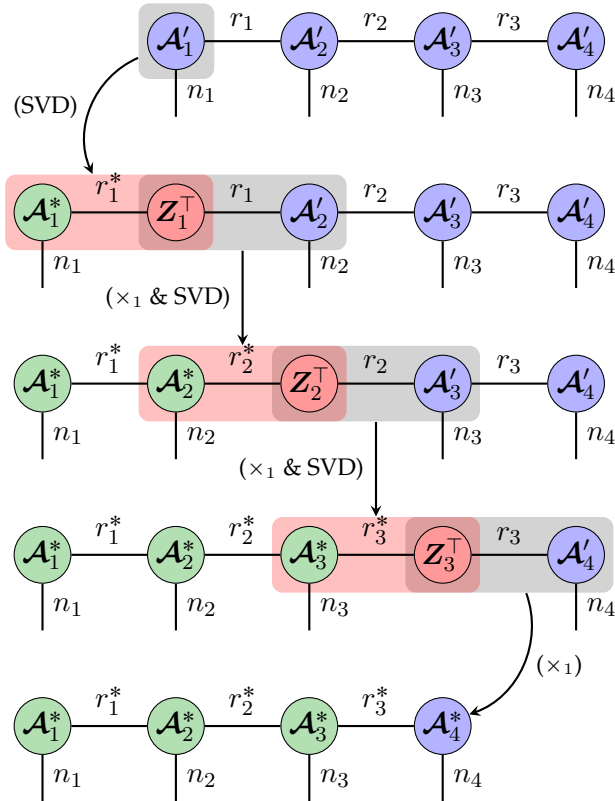


FIGURE 4.3: Tensor network diagram for the compression of $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ using Algorithm 7. Green nodes have second unfoldings with orthonormal columns, whereas blue nodes have first unfoldings with orthonormal rows.

We show in Figure 4.3 how the procedure is carried out on a $n_1 \times \dots \times n_d$ right-orthogonal tensor \mathcal{A}' with TT-ranks $(1, r_1, r_2, r_3, 1)$. In particular, we can observe that each of the “local” right factors \mathcal{Z}_j (for $j \in \{1, \dots, d-1\}$) acts as an interface between a left-orthogonal and a right-orthogonal sub-tensor, represented respectively in green and blue.

4.2.3 Rounding of tensors in the TT format in two sweeps

Algorithm 7 can be easily modified to tackle the rounding problem: in this case the final TT-ranks are unknown, but we want to guarantee that the relative error (4.10) is smaller than some tolerance $\varepsilon > 0$. To this aim, we compute each new rank r_j^* (for $j \in \{1, \dots, d-1\}$) in such a way that the approximation has the adjusted relative accuracy δ (with $\delta \leq \varepsilon$): if we proceed from left to right, this can be expressed as

$$\frac{\|U_j(:, 1:r_j^*)\Sigma_j(1:r_j^*, 1:r_j^*)V_j(:, 1:r_j^*)^\top - (\mathcal{A}_j^{**})^{\langle 2 \rangle}\|_F}{\|\mathcal{A}\|_F} \leq \delta \quad (4.17)$$

Theorem 4.3 below shows that δ is within a factor (which depends on the order d) of the required accuracy ε . Accordingly, to solve the rounding problem with precision ε , it is enough to run Algorithm 7 with condition (4.17) and $\delta = \varepsilon/\sqrt{d-1}$ (see Algorithm 8).

Algorithm 8: Rounding of a tensor in the TT format

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, $\varepsilon > 0$

Result: left-orthogonal $\mathcal{A}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1^*, \dots, \mathcal{A}_d^*$ such that the relative error (4.10) is at most ε

```

1  $\delta := \varepsilon/\sqrt{d-1}$ 
2  $\mathcal{A}' := \text{Algorithm 6}(\mathcal{A}, \text{'right'})$ 
3  $\mathcal{A}_1^{**} := \mathcal{A}'_1$ 
4 for  $j = 1, 2, \dots, d-1$  do
5   Compute SVD:  $(\mathcal{A}_j^{**})^{\langle 2 \rangle} =: U_j \Sigma_j V_j^\top$ 
6   Find the smallest  $r_j^*$  such that (4.17) holds
7    $\mathcal{W}_j := U_j(:, 1:r_j^*)$ 
8    $\mathcal{Z}_j := V_j(:, 1:r_j^*) \Sigma_j(1:r_j^*, 1:r_j^*)$ 
9    $\mathcal{A}_j^* := \text{reshape}(\mathcal{W}_j, [r_{j-1}^*, n_j, r_j^*])$ 
10   $\mathcal{A}_{j+1}^{**} := \mathcal{A}'_{j+1} \times_1 \mathcal{Z}_j^\top$ 
11 end
12  $\mathcal{A}_d^* := \mathcal{A}_d^{**}$ 

```

Theorem 4.3. *The tensor \mathcal{A}^* resulting from Algorithm 8 is such that*

$$\|\mathcal{A}^* - \mathcal{A}\|_F \leq \varepsilon \|\mathcal{A}\|_F \quad (4.18)$$

Proof. We follow the idea behind the proof of Theorem 2.2 in [Ose11], and we employ the same notation used in Algorithm 8.

For each $j \in \{0, \dots, d-1\}$, we define the partial approximation $\mathcal{A}^{(j)}$ as the $n_1 \times \dots \times n_d$ tensor which admits the TT decomposition with cores

$$\mathcal{A}_1^*, \dots, \mathcal{A}_j^*, \mathcal{A}_{j+1}^{**}, \mathcal{A}'_{j+2}, \dots, \mathcal{A}'_d$$

where

$$\mathcal{A}_{j+1}^{**} = \mathcal{A}'_{j+1} \times_1 \mathbf{Z}_j^\top \quad \text{for } j = 1, \dots, d-1$$

and $\mathcal{A}_1^{**} = \mathcal{A}'_1$.

It is important to observe that, due to the properties of 1st and 3rd mode products, $\mathcal{A}^{(j)}$ also admits the alternative TT decomposition with cores

$$\mathcal{A}_1^*, \dots, \mathcal{A}_{j-1}^*, \mathcal{A}_j^* \times_3 \mathbf{Z}_j, \mathcal{A}'_{j+1}, \dots, \mathcal{A}'_d$$

for any $j \in \{1, \dots, d-1\}$.

With this definition, each $\mathcal{A}^{(j)}$ corresponds to the approximation of \mathcal{A} immediately after the j^{th} iteration of Algorithm 8. In particular, $\mathcal{A}^{(0)} = \mathcal{A}$ and $\mathcal{A}^{(d-1)} = \mathcal{A}^*$.

Since a tensor and all of its unfoldings have the same Frobenius norm, we can express the error as

$$\|\mathcal{A}^* - \mathcal{A}\|_F = \|\mathcal{A}^{(d-1)} - \mathcal{A}^{(0)}\|_F = \left\| \left(\mathcal{A}^{(d-1)} - \mathcal{A}^{(1)} \right)^{\langle 1 \rangle} + \left(\mathcal{A}^{(1)} - \mathcal{A}^{(0)} \right)^{\langle 1 \rangle} \right\|_F$$

The two terms can be manipulated as

$$\left(\mathcal{A}^{(d-1)} - \mathcal{A}^{(1)} \right)^{\langle 1 \rangle} = (\mathcal{A}_1^*)^{\langle 2 \rangle} \left(\mathcal{A}_{\geq 2}^{(d-1)} - \mathcal{A}_{\geq 2}^{(1)} \right)^\top = \mathbf{W}_1 \left(\mathcal{A}_{\geq 2}^{(d-1)} - \mathcal{A}_{\geq 2}^{(1)} \right)^\top$$

and

$$\begin{aligned} \left(\mathcal{A}^{(1)} - \mathcal{A}^{(0)} \right)^{\langle 1 \rangle} &= (\mathcal{A}_1^* \times_3 \mathbf{Z}_1 - \mathcal{A}_1^{**})^{\langle 2 \rangle} \mathcal{A}'_{\geq 2}{}^\top \\ &= \left(\mathbf{W}_1 \mathbf{Z}_1^\top - (\mathcal{A}_1^{**})^{\langle 2 \rangle} \right) \mathcal{A}'_{\geq 2}{}^\top \\ &= \left(\mathbf{W}_1 \mathbf{W}_1^\top (\mathcal{A}_1^{**})^{\langle 2 \rangle} - (\mathcal{A}_1^{**})^{\langle 2 \rangle} \right) \mathcal{A}'_{\geq 2}{}^\top \\ &= (\mathbf{W}_1 \mathbf{W}_1^\top - \mathbf{I}_{n_1}) (\mathcal{A}_1^{**})^{\langle 2 \rangle} \mathcal{A}'_{\geq 2}{}^\top \end{aligned}$$

Since \mathbf{W}_1 has orthonormal columns, it holds that

$$\mathbf{W}_1^\top (\mathbf{W}_1 \mathbf{W}_1^\top - \mathbf{I}_{n_1}) = \mathbf{0}$$

Hence, the Pythagorean theorem can be applied:

$$\begin{aligned} \|\mathcal{A}^* - \mathcal{A}\|_F^2 &= \left\| \left(\mathcal{A}^{(d-1)} - \mathcal{A}^{(1)} \right)^{\langle 1 \rangle} \right\|_F^2 + \left\| \left(\mathcal{A}^{(1)} - \mathcal{A}^{(0)} \right)^{\langle 1 \rangle} \right\|_F^2 \\ &= \left\| \mathcal{A}^{(d-1)} - \mathcal{A}^{(1)} \right\|_F^2 + \left\| \mathbf{W}_1 \mathbf{Z}_1^\top - (\mathcal{A}_1^{**})^{\langle 2 \rangle} \right\|_F^2 \end{aligned}$$

Now we consider the second unfolding of $\mathcal{A}^{(d-1)} - \mathcal{A}^{(1)}$ and proceed similarly, obtaining

$$\left(\mathcal{A}^{(d-1)} - \mathcal{A}^{(2)} \right)^{\langle 2 \rangle} = (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) \mathbf{W}_2 \left(\mathcal{A}_{\geq 3}^{(d-1)} - \mathcal{A}_{\geq 3}^{(2)} \right)^\top$$

and

$$\left(\mathcal{A}^{(2)} - \mathcal{A}^{(1)}\right)^{\langle 2 \rangle} = (\mathbf{I}_{n_2} \otimes \mathbf{W}_1) \left(\mathbf{W}_2 \mathbf{W}_2^\top - \mathbf{I}_{n_2 r_1^*} \right) (\mathcal{A}_2^{**})^{\langle 2 \rangle} \mathcal{A}'_{\geq 3}{}^\top$$

After removing the common term $\mathbf{I}_{n_2} \otimes \mathbf{W}_1$, which has orthonormal columns and does not affect the Frobenius norm, the remaining leading terms are orthogonal, and the Pythagorean theorem yields

$$\|\mathcal{A}^{(d-1)} - \mathcal{A}^{(1)}\|_F^2 = \|\mathcal{A}^{(d-1)} - \mathcal{A}^{(2)}\|_F^2 + \|\mathbf{W}_2 \mathbf{Z}_2^\top - (\mathcal{A}_2^{**})^{\langle 2 \rangle}\|_F^2$$

A recursive argument leads to the expression

$$\|\mathcal{A}^* - \mathcal{A}\|_F^2 = \sum_{j=1}^{d-1} \|\mathbf{W}_j \mathbf{Z}_j^\top - (\mathcal{A}_j^{**})^{\langle 2 \rangle}\|_F^2$$

from which, by using the uniform bound (4.17), we derive

$$\|\mathcal{A}^* - \mathcal{A}\|_F^2 \leq (d-1)\delta^2 \|\mathcal{A}\|_F^2 = \varepsilon^2 \|\mathcal{A}\|_F^2$$

□

4.3 Compression of tensors in the TT format using a randomized approach

In the last section we have described the *state-of-the-art* deterministic algorithms to compress tensors in the TT format, and we have shown that the accuracy and efficiency of the procedure rely on the first “sweep”, which is used to orthogonalize the tensor. In the same spirit as the one applied to matrices in Chapter 3, we may hope to reduce the complexity of the algorithm from $\mathcal{O}(dnr^3)$ to $\mathcal{O}(dnr^2r^*)$. However, the intrinsic cost ($\mathcal{O}(dnr^3)$ flops) due to the orthogonalization step cannot be avoided in a deterministic approach.

In this section we describe a randomized algorithm which achieves the desired complexity $\mathcal{O}(dnr^2r^*)$, by exploiting the stunning efficiency with which some operations can be executed in the TT format.

4.3.1 Multidimensional tensor contraction

First, we introduce the operation which the algorithm relies on.

Definition 4.3. Let \mathcal{A} be a tensor of size $n_1 \times \dots \times n_d$, and let $j \in \{0, \dots, d-1\}$. Consider $d-j$ arbitrary vectors $\{\mathbf{u}_l\}_{l=j+1}^d$, with $\mathbf{u}_l \in \mathbb{R}^{n_l}$ for $l \in \{j+1, \dots, d\}$.

We define the mode- $(j+1, \dots, d)$ contraction of \mathcal{A} onto the rank-1 tensor $\bigotimes_{l=j+1}^d \mathbf{u}_l$ as the $n_1 \times \dots \times n_j$ tensor (or a scalar if $j=0$) with entries defined by

$$\left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)}(i_1, \dots, i_j) = \sum_{i_{j+1}, \dots, i_d} \mathcal{A}(i_1, \dots, i_d) \mathbf{u}_{j+1}(i_{j+1}) \dots \mathbf{u}_d(i_d) \quad (4.19)$$

for $(i_1, \dots, i_j) \in I_1 \times \dots \times I_j$, with the sum over i_l being carried out between 1 and n_l , for $l \in \{j+1, \dots, d\}$.

It is trivial to see that the contraction $\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \rangle_{(j+1, \dots, d)}$ can be obtained by removing all the singleton dimensions from

$$(\mathcal{A} \times_{j+1} \mathbf{u}_{j+1}^\top \times_{j+2} \mathbf{u}_{j+2}^\top \dots \times_d \mathbf{u}_d^\top) \in \mathbb{R}^{n_1 \times \dots \times n_j \times 1 \times \dots \times 1}$$

In particular, if the tensor \mathcal{A} is in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, the terms appearing in the sum (4.19) are separable:

$$\begin{aligned} & \left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)}(i_1, \dots, i_j) = \\ & = \mathcal{A}_1(i_1) \dots \mathcal{A}_j(i_j) \left(\sum_{i_{j+1}} \mathcal{A}_{j+1}(i_{j+1}) \mathbf{u}_{j+1}(i_{j+1}) \right) \dots \left(\sum_{i_d} \mathcal{A}_d(i_d) \mathbf{u}_d(i_d) \right) \end{aligned} \quad (4.20)$$

From this expression, we can devise a very efficient algorithm to compute contractions in the TT format, whose original version can be found in [Ose11], where only full contractions are considered.

We proceed in this way: let \mathcal{A} be a tensor of size $n_1 \times \dots \times n_d$, whose TT decomposition has cores $\mathcal{A}_1, \dots, \mathcal{A}_d$ and TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$. Also, let $\{\mathbf{u}_l\}_{l=j+1}^d$ be such that $\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \rangle_{(j+1, \dots, d)}$ is defined.

Following (4.20), we compute the vector

$$\mathbf{v}_d = \sum_{i_d} \mathcal{A}_d(i_d) \mathbf{u}_d(i_d) \in \mathbb{R}^{r_{d-1}}$$

and we proceed recursively from right to left

$$\mathbf{v}_l = \sum_{i_l} (\mathcal{A}_l(i_l) \mathbf{v}_{l+1}) \mathbf{u}_l(i_l) \in \mathbb{R}^{r_{l-1}} \quad \text{for } l = d-1, \dots, j+1$$

Now it suffices to update the j^{th} core

$$\mathcal{A}'_j = \mathcal{A}_j \times_3 \mathbf{v}_{j+1}^\top \in \mathbb{R}^{r_{j-1} \times n_j \times 1}$$

so that the TT decomposition of the contraction is

$$\left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)}(i_1, \dots, i_j) = \mathcal{A}_1(i_1) \dots \mathcal{A}_{j-1}(i_{j-1}) \mathcal{A}'_j(i_j)$$

This procedure is summarized in Algorithm 9, and costs $\mathcal{O}((d-j)nr^2)$ flops.

This algorithm also admits an intriguing graphical representation using tensor network diagrams. In Figure 4.4 we consider a tensor with size $n_1 \times \dots \times n_5$ and TT-ranks $(1, r_1, \dots, r_4, 1)$, which is contracted along the modes (3, 4, 5) onto the rank-1 tensor $\bigotimes_{l=3}^5 \mathbf{u}_l$: in the first step, we collapse \mathcal{A}_5 onto \mathbf{u}_5 to obtain \mathbf{v}_5 ; then $\mathcal{A}_4, \mathbf{u}_4$ and \mathbf{v}_5 are used to compute \mathbf{v}_4 and so on. At the end of the procedure, we collapse the second core onto \mathbf{v}_3 to obtain the desired TT representation.

Algorithm 9: Mode- $(j + 1, \dots, d)$ contraction of a tensor in the TT format

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$ and TT-ranks $(1, r_1, \dots, r_{d-1}, 1)$, index $j \in \{0, \dots, d-1\}$, vectors $\{\mathbf{u}_l\}_{l=j+1}^d$ of suitable size

Result: $\mathcal{A}' = \langle \mathcal{A}, \otimes_{l=j+1}^d \mathbf{u}_l \rangle_{(j+1, \dots, d)} \in \mathbb{R}^{n_1 \times \dots \times n_j}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_{j-1}, \mathcal{A}'_j$

```

1  $v_{d+1} := 1$ 
2 for  $l = d, d-1, \dots, j+1$  do
3    $\mathbf{v}_l := \text{zeros}(r_{l-1}, 1)$ 
4   for  $i_l = 1, 2, \dots, n_l$  do
5      $\mathbf{v}_l := \mathbf{v}_l + (\mathcal{A}_l(i_l) \mathbf{v}_{l+1}) \mathbf{u}_l(i_l)$ 
6   end
7 end
8  $\mathcal{A}'_j := \mathcal{A}_j \times_3 \mathbf{v}_{j+1}^\top$ 
    
```

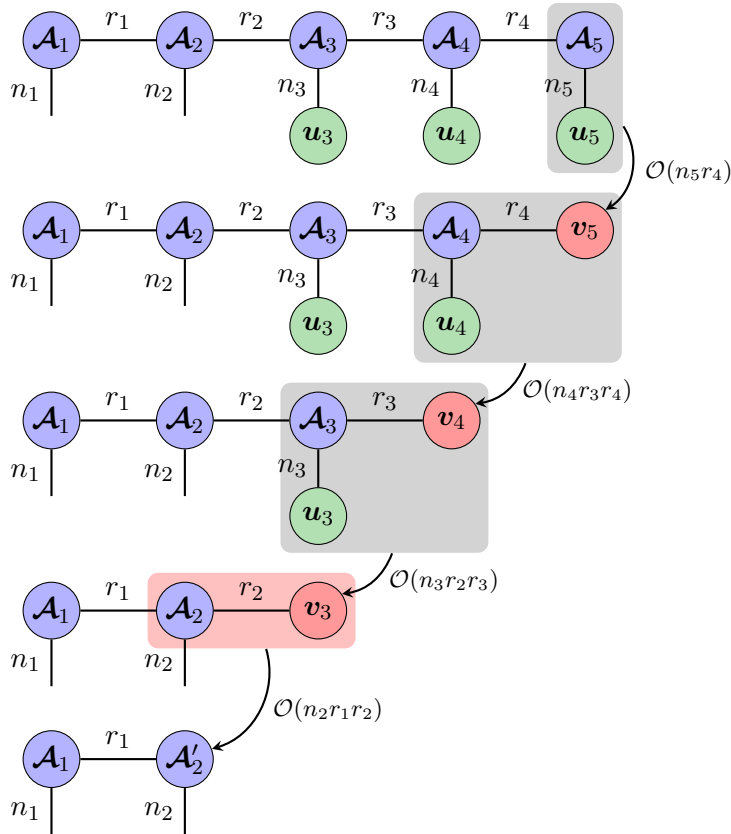


FIGURE 4.4: Tensor network diagram for the contraction of $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_5}$ along the modes $(3, 4, 5)$ onto the rank-1 tensor $\otimes_{l=3}^5 \mathbf{u}_l$. At each step, one or two edges are collapsed according to Algorithm 9.

In particular, we may easily evaluate the complexity of the algorithm from the diagram, as follows. For each value of $l \in \{d, d-1, \dots, j+1\}$ we identify the set N_l of nodes (either cores of the tensor, vectors \mathbf{u}_l , or partial results of the procedure) which are used in the corresponding step of Algorithm 9: such sets are represented as gray rectangles in Figure 4.4. With this definition, the complexity of the step can be simply obtained as the product of the labels of all the edges having (at least) an element of N_l as an endpoint.

Among the many applications of multidimensional contractions, we consider one which is very useful in the context of tensor compression (see Section 4.3.2): given a tensor \mathcal{A} , contractions can be applied to compute the matrix-vector product between any unfolding $\mathcal{A}^{<j>}$ and structured “rank-1” vectors, as stated by the following result.

Lemma 4.4. *Let \mathcal{A} , j and $\{\mathbf{u}_l\}_{l=j+1}^d$ be as in Definition 4.3. Let the vector $\mathbf{u}^* \in \mathbb{R}^{n_{j+1} \dots n_d}$ be the Kronecker product*

$$\mathbf{u}^* = \mathbf{u}_d \otimes \dots \otimes \mathbf{u}_{j+1} \quad (4.21)$$

Then

$$\mathcal{A}^{<j>} \mathbf{u}^* = \left(\left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)} \right)^{<j>} \quad (4.22)$$

Proof. Let i_{row} and i_{col} be defined as in Definition 2.1. Since i_{col} and the Kronecker product (4.21) induce, respectively, a colexicographical and a lexicographical ordering of the indices (i_{j+1}, \dots, i_d) , we have that

$$\mathbf{u}^*(i_{col}(i_{j+1}, \dots, i_d)) = \mathbf{u}_{j+1}(i_{j+1}) \dots \mathbf{u}_d(i_d) \quad \text{for } (i_{j+1}, \dots, i_d) \in I_{j+1} \times \dots \times I_d$$

Hence

$$\begin{aligned} (\mathcal{A}^{<j>} \mathbf{u}^*)(i_{row}(i_1, \dots, i_j)) &= \sum_{\alpha=1}^{n_{j+1} \dots n_d} \mathcal{A}^{<j>}(i_{row}(i_1, \dots, i_j), \alpha) \mathbf{u}^*(\alpha) \\ &= \sum_{i_{j+1}, \dots, i_d} \mathcal{A}(i_1, \dots, i_d) \mathbf{u}_{j+1}(i_{j+1}) \dots \mathbf{u}_d(i_d) \\ &= \left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)}(i_1, \dots, i_j) \\ &= \left(\left\langle \mathcal{A}, \bigotimes_{l=j+1}^d \mathbf{u}_l \right\rangle_{(j+1, \dots, d)} \right)^{<j>}(i_{row}(i_1, \dots, i_j)) \end{aligned}$$

for $(i_1, \dots, i_j) \in I_1 \times \dots \times I_j$. □

Moreover, it is also important to fully understand the meaning of the partial results \mathbf{v}_k for $k \in \{j+2, \dots, d\}$ in Algorithm 9. It is easy to check that they appear in the expansion of contractions along fewer modes: indeed, it is trivial to verify that, for any $k \in \{j, \dots, d-1\}$, it holds

$$\left\langle \mathcal{A}, \bigotimes_{l=k+1}^d \mathbf{u}_l \right\rangle_{(k+1, \dots, d)}(i_1, \dots, i_k) = \mathcal{A}_1(i_1) \dots \mathcal{A}_{k-1}(i_{k-1}) \mathcal{A}_k(i_k) \mathbf{v}_{k+1} \quad (4.23)$$

for $(i_1, \dots, i_k) \in I_1 \times \dots \times I_k$.

4.3.2 Truncation of tensors in the TT format in a single sweep

Now that we have an efficient way to compute the action of any unfolding onto a structured vector, we can try to apply the ideas introduced in Section 3.1 to compress any unfolding of the target tensor.

Given the usual tensor \mathcal{A} in the TT format (4.2), we consider the first unfolding $\mathcal{A}^{<1>} = \mathcal{A}_{\leq 1} \mathcal{A}_{\geq 2}^\top$. We seek a rank- r_1^* dyadic approximation of the form

$$\mathcal{A}^{<1>} \approx \mathbf{W}_1 \mathbf{Z}_1^\top \quad \text{with } \mathbf{W}_1^\top \mathbf{W}_1 = \mathbf{I}_{r_1^*}$$

This time, we cannot identify the left factor \mathbf{W}_1 just from the left interface $\mathcal{A}_{\leq 1}$, since in general the TT decomposition of \mathcal{A} is not right-orthogonal: instead, we need to use the whole unfolding (in this issue resides the reason why the orthogonalization step is required to guarantee efficiency in the deterministic approach).

We do this in a randomized fashion, as follows. We multiply the matrix by $r_1^* + p$ random vectors (with $p > 0$ being an oversampling parameter):

$$\mathbf{y}^{(l)} = \mathcal{A}^{<1>} \boldsymbol{\omega}^{(l)} \quad \text{for } l \in \{1, \dots, r_1^* + p\}$$

From the resulting samples we derive an approximation for the left factor \mathbf{W}_1 , e.g. by computing the orthogonal factor of the QR decomposition of the matrix $[\mathbf{y}^{(1)} | \dots | \mathbf{y}^{(r_1^* + p)}]$.

If each random vector $\boldsymbol{\omega}^{(l)} \in \mathbb{R}^{n_2 \dots n_d}$ were drawn from a Gaussian distribution, computing the corresponding sample would require $\mathcal{O}(rn^{d-1})$ flops, making the algorithm extremely inefficient. Instead, in the wake of Lemma 4.4, we choose random vectors with a “rank-1” structure: to generate each sample $\boldsymbol{\omega}^{(l)}$, for $l \in \{1, \dots, r_1^* + p\}$, we draw $d - 1$ Gaussian random vectors $\{\boldsymbol{\omega}_j^{(l)}\}_{j=2}^d$ (with $\boldsymbol{\omega}_j^{(l)} \in \mathbb{R}^{n_j}$ for $j \in \{2, \dots, d\}$), and we take

$$\boldsymbol{\omega}^{(l)} = \boldsymbol{\omega}_d^{(l)} \otimes \dots \otimes \boldsymbol{\omega}_2^{(l)}$$

With this definition, the computation of each of the samples $\{\mathbf{y}^{(l)}\}_{l=1}^{r_1^* + p}$ can be carried out very efficiently as a multidimensional contraction, thanks to Lemma 4.4.

Once the samples have been computed, they are orthogonalized and arranged in a $n_1 \times (r_1^* + p)$ matrix \mathbf{W}'_1 with orthonormal columns, so that we have the approximate rank- $(r_1^* + p)$ decomposition

$$\mathcal{A}^{<1>} \approx \mathbf{W}'_1 \left(\mathbf{W}'_1{}^\top \mathcal{A}^{<1>} \right) = \mathbf{W}'_1 \left(\mathbf{W}'_1{}^\top (\mathcal{A}_1)^{<2>} \mathcal{A}_{\geq 2}^\top \right)$$

Similarly to Algorithms 6 and 7, \mathbf{W}'_1 is used to replace the first core of the tensor, so that $(\mathcal{A}'_1)^{<2>} = \mathbf{W}'_1$, and we use the right factor of the rank- $(r_1^* + p)$ approximation of $(\mathcal{A}_1)^{<2>}$, i.e.

$$\mathbf{Z}'_1 = (\mathcal{A}_1)^{<2>}{}^\top \mathbf{W}'_1,$$

to update the second core:

$$\mathcal{A}''_2 = \mathcal{A}_2 \times_1 \mathbf{Z}'_1{}^\top$$

Now we consider the second unfolding

$$\mathcal{A}^{<2>} = \mathcal{A}_{\leq 2} \mathcal{A}_{\geq 3}^\top = (\mathbf{I}_{n_2} \otimes \mathcal{A}_{\leq 1}) (\mathcal{A}_2)^{<2>} \mathcal{A}_{\geq 3}^\top$$

$$\begin{aligned}
&\approx \left(\mathbf{I}_{n_2} \otimes \left(\mathbf{W}'_1 \mathbf{Z}'_1{}^\top \right) \right) (\mathcal{A}_2)^{\langle 2 \rangle} \mathcal{A}_{\geq 3}^\top \\
&= \left(\mathbf{I}_{n_2} \otimes \mathbf{W}'_1 \right) \left(\mathbf{I}_{n_2} \otimes \mathbf{Z}'_1{}^\top \right) (\mathcal{A}_2)^{\langle 2 \rangle} \mathcal{A}_{\geq 3}^\top \\
&= \left(\mathbf{I}_{n_2} \otimes \mathbf{W}'_1 \right) (\mathcal{A}_2'')^{\langle 2 \rangle} \mathcal{A}_{\geq 3}^\top
\end{aligned}$$

Since \mathbf{W}'_1 has orthonormal columns, so does $\mathbf{I}_{n_2} \otimes \mathbf{W}'_1$, and it suffices to estimate the range of $(\mathcal{A}_2'')^{\langle 2 \rangle} \mathcal{A}_{\geq 3}^\top$. Once again, this can be done using multidimensional contractions: indeed, it is sufficient to start the computation of the new samples

$$\mathbf{y}^{(l)} = \mathcal{A}^{\langle 2 \rangle} \left(\boldsymbol{\omega}_d^{(l)} \otimes \dots \otimes \boldsymbol{\omega}_3^{(l)} \right) \quad \text{for } l \in \{1, \dots, r_2^* + p\}$$

with $\{\boldsymbol{\omega}_j^{(l)}\}_{j=3}^d$ (for $l \in \{1, \dots, r_2^* + p\}$) being a suitable subset of the random vectors used in the previous step.

In particular, we are not interested in the contributions given by the original first and second cores: thus, we just compute the value of v_3 in Algorithm 9, which represents one of the partial samples

$$\bar{\mathbf{y}}^{(l)} = \mathcal{A}_{\geq 3}^\top \left(\boldsymbol{\omega}_d^{(l)} \otimes \dots \otimes \boldsymbol{\omega}_3^{(l)} \right) \quad \text{for } l \in \{1, \dots, r_2^* + p\}$$

Then we multiply each of these vectors (from the left) by $(\mathcal{A}_2'')^{\langle 2 \rangle}$ to get the desired results. The final samples are orthonormalized and arranged in a $(n_2(r_1^* + p)) \times (r_2^* + p)$ matrix \mathbf{W}'_2 (with orthonormal columns), so that we have the approximate rank- $(r_2^* + p)$ decomposition

$$\mathcal{A}^{\langle 2 \rangle} \approx \left((\mathbf{I}_{n_2} \otimes \mathbf{W}'_1) \mathbf{W}'_2 \right) \left(\mathbf{W}'_2{}^\top (\mathcal{A}_2'')^{\langle 2 \rangle} \mathcal{A}_{\geq 3}^\top \right)$$

Now the second core is overwritten in such a way that $(\mathcal{A}_2')^{\langle 2 \rangle} = \mathbf{W}'_2$, and the third core is updated using

$$\mathbf{Z}'_2 = (\mathcal{A}_2'')^{\langle 2 \rangle \top} \mathbf{W}'_2$$

so that

$$\mathcal{A}_3'' = \mathcal{A}_3 \times_1 \mathbf{Z}'_2{}^\top$$

Then this process is repeated on all the cores.

Since this results in a tensor with TT-ranks $(1, r_1^* + p, \dots, r_{d-1}^* + p, 1)$, we carry out a (deterministic) truncation of the tensor down to the desired TT-ranks. In particular, since the tensor \mathcal{A}' is left-orthogonal by construction, it suffices to apply Algorithm 7 in reverse, without the need for a preliminary orthogonalization step.

The resulting procedure is summarized in Algorithm 11, where the subroutine described in Algorithm 10 is used to compute the random samples through multidimensional contractions. Since Algorithm 10 computes the samples recursively, its cost is $\mathcal{O}(dnr^2k)$ flops, which results in $\mathcal{O}(dnr^2(r^* + p))$ flops when applied within Algorithm 11.

This operation, and the update in line 10 (which costs $\mathcal{O}(nr^2(r^* + p))$ flops per iteration), are responsible for the overall complexity of the algorithm, which equals $\mathcal{O}(dnr^2(r^* + p))$.

Algorithm 10: Generation of samples for the range of tensor unfoldings

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$ and TT-ranks $(r_0 = 1, r_1, \dots, r_{d-1}, r_d = 1)$, number of samples k

Result: $\{\mathbf{Y}_j\}_{j=2}^d$, with $\mathbf{Y}_j = \mathcal{A}_{\geq j}^\top \Omega_j \in \mathbb{R}^{r_{j-1} \times k}$ for a structured random matrix $\Omega_j \in \mathbb{R}^{(n_j \dots n_d) \times k}$ (for $j = 2, \dots, d$)

- 1 $\mathbf{Y}_{d+1} := [1 | \dots | 1]$ % of size $1 \times k$
- 2 **for** $j = d, d-1, \dots, 2$ **do**
- 3 Generate Gaussian matrix $\Omega_j = [\omega_j^{(1)} | \dots | \omega_j^{(k)}] \in \mathbb{R}^{n_j \times k}$
- 4 $\mathbf{Y}_j := \text{zeros}(r_{j-1}, k)$
- 5 **for** $l = 1, \dots, k$ **do**
- 6 **for** $i_j = 1, \dots, n_j$ **do**
- 7 $\mathbf{Y}_j(:, l) := \mathbf{Y}_j(:, l) + (\mathcal{A}_j(i_j) \mathbf{Y}_{j+1}(:, l)) \omega_j^{(l)}(i_j)$
- 8 **end**
- 9 **end**
- 10 **end**

Algorithm 11: Randomized truncation of a tensor in the TT format

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, desired TT-ranks $(r_0^* = 1, r_1^*, \dots, r_{d-1}^*, r_d^* = 1)$, integer $p \geq 0$

Result: $\mathcal{A}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1^*, \dots, \mathcal{A}_d^*$ with the specified TT-ranks, such that $\mathcal{A} \approx \mathcal{A}^*$

- 1 $k := \max\{r_1^*, \dots, r_{d-1}^*\} + p$
- 2 $\{\mathbf{Y}_2, \dots, \mathbf{Y}_d\} := \text{Algorithm 10}(\mathcal{A}, k)$
- 3 $\mathcal{A}_1'' := \mathcal{A}_1$
- 4 $r_{old}^* := 1$
- 5 **for** $j = 1, 2, \dots, d-1$ **do**
- 6 $\mathbf{Y}_j' := (\mathcal{A}_j'')^{<2>} \mathbf{Y}_{j+1}(:, r_j^* + p)$
- 7 Compute QR decomposition: $\mathbf{Y}_j' =: \mathbf{W}_j' \mathbf{R}_j$
- 8 $\mathcal{A}_j' := \text{reshape}(\mathbf{W}_j', [r_{old}^*, n_j, r_j^* + p])$
- 9 $\mathbf{Z}_j' := (\mathcal{A}_j'')^{<2>^\top} \mathbf{W}_j'$
- 10 $\mathcal{A}_{j+1}'' := \mathcal{A}_{j+1} \times_1 \mathbf{Z}_j'^\top$
- 11 $r_{old}^* := r_j^* + p$
- 12 **end**
- 13 $\mathcal{A}_d' := \mathcal{A}_d''$
- 14 $\mathcal{A}^* := \text{Algorithm 7}(\mathcal{A}', (1, r_1^*, \dots, r_{d-1}^*, 1))$ % from right to left

4.3.3 Randomized rounding of tensors in the TT format

By applying Algorithm 11, it is possible to truncate a tensor down to some prescribed TT-ranks. In this section we describe a possible way to extend this procedure to tackle the rounding problem with precision $\varepsilon > 0$.

Given the usual tensor \mathcal{A} in the TT-format (4.2), we wish to find a tensor \mathcal{A}^* with smaller TT-ranks, such that the relative Frobenius error (4.10) is smaller than ε . Due to the trivial generalization of Theorem 4.3 to non-orthogonal tensors, it is sufficient to require that

$$\frac{\|\mathbf{W}'_j \mathbf{W}'_j{}^\top (\mathcal{A}''_j)^{\langle 2 \rangle} \mathcal{A}^\top_{\geq j+1} - (\mathcal{A}''_j)^{\langle 2 \rangle} \mathcal{A}^\top_{\geq j+1}\|_F}{\|\mathcal{A}\|_F} \leq \frac{\varepsilon}{\sqrt{d-1}} \quad (4.24)$$

for $j \in \{1, \dots, d-1\}$, using the same notation as in Algorithm 11.

Unfortunately, in this case it is quite expensive to compute the left hand side, or even to estimate it with reasonable accuracy, due to the presence of the right interface matrix $\mathcal{A}_{\geq j+1}$. In this section we describe an alternative adaptive approach, which achieves a cost very similar to the truncation case.

Let $\delta = \varepsilon/\sqrt{d-1}$, and assume that \hat{r}_j is the δ -rank of $\mathcal{A}^{\langle j \rangle}$ (for $j \in \{1, \dots, d-1\}$), i.e. it is the smallest integer such that there exists a $(n_1 \cdots n_j) \times \hat{r}_j$ matrix \mathbf{U}_j with orthonormal columns such that

$$\frac{\|(\mathbf{U}_j \mathbf{U}_j{}^\top - \mathbf{I}_{n_1 \cdots n_j}) \mathcal{A}^{\langle j \rangle}\|_F}{\|\mathcal{A}\|_F} \leq \delta$$

Theorems 4.1 and 4.3 ensure that there exists a tensor with TT-ranks $(1, \hat{r}_1, \dots, \hat{r}_{d-1}, 1)$ which solves the rounding problem.

Consider a vector $(1, r_1^{(0)}, \dots, r_{d-1}^{(0)}, 1)$ containing guesses for the TT-ranks of the compressed tensor. We start by approximating \mathcal{A} with a tensor \mathcal{A}' having the specified TT-ranks. To this aim, we run Algorithm 11 with $p = 0$ (in particular we skip the recompression step, i.e. line 14).

Now we wish to know whether \mathcal{A}' achieves (4.24) for $j \in \{1, \dots, d-1\}$. We may start by checking if the condition

$$r_j^{(0)} \geq \hat{r}_j \quad \text{for } j \in \{1, \dots, d-1\}$$

holds. However, the ranks $\{\hat{r}_j\}_{j=1}^{d-1}$ are unknown (and it would be expensive to compute them). Hence, we follow a different route: using Algorithm 8, we recompress \mathcal{A}' down to precision ε , and check if all the TT-ranks (except r_0 and r_d) decrease by at least $q > 0$.

Here we are assuming that, in order to be a good (enough) approximation of \mathcal{A} , the tensor \mathcal{A}' should be able to identify accurately the first $\hat{r}_j + q$ singular values of $\mathcal{A}^{\langle j \rangle}$, for $j \in \{1, \dots, d-1\}$. In particular, the parameter q is used to counter two effects due to the randomness of the algorithm: foremost, the uncertainty in the approximation of the first \hat{r}_j singular values (and vectors) of the unfolding $\mathcal{A}^{\langle j \rangle}$, since they represent the true objective of the compression; in second measure, the error in the approximation of the remaining $r_j^{(0)} - \hat{r}_j$ singular values, which may

cause the recompression to stop earlier than required, due to a poor estimation of $\|\mathcal{A}' - \mathcal{A}\|_F$ and $\|\mathcal{A}\|_F$.

As such, the value of q should be prescribed based on the available information on the spectra of the unfoldings. For exponentially decaying singular values, a small value of q is usually enough (see Section 4.4.3), whereas a larger value should be used for slower decays.

If the check is successful, we accept \mathcal{A}' (or its recompressed version) as the solution of the rounding problem. Otherwise, we increase all the dubious ranks (i.e. those which do not decrease enough through the recompression step) by some fixed amount $m > 0$. This results in a vector containing the new guesses $(1, r_1^{(1)}, \dots, r_{d-1}^{(1)}, 1)$.

Now we restart the procedure using the updated guesses for the TT-ranks, and we repeat this truncation–rounding–update cycle until the condition on the TT-ranks is satisfied (see Algorithm 12).

Algorithm 12: Randomized rounding of a tensor in the TT format

Data: tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$, $\varepsilon > 0$,
 guess for the TT-ranks $(1, r_1^{(0)}, \dots, r_{d-1}^{(0)}, 1)$, positive integers q and m
Result: $\mathcal{A}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1^*, \dots, \mathcal{A}_d^*$ such that the
 relative error (4.10) is at most ε with high probability

```

1  $l := 0$ 
2 do
3    $\mathcal{A}' := \text{Algorithm 11}(\mathcal{A}, (1, r_1^{(l)}, \dots, r_{d-1}^{(l)}, 1), 0)$ 
4    $\mathcal{A}^* := \text{Algorithm 8}(\mathcal{A}', \varepsilon)$  % from right to left
5    $(1, r_1^*, \dots, r_{d-1}^*, 1) = \text{TT-ranks}(\mathcal{A}^*)$ 
6    $\text{changed} := \text{false}$ 
7   for  $j = 1, \dots, d - 1$  do
8     if  $r_j^{(l)} \geq r_j^* + q$  then
9        $r_j^{(l+1)} := r_j^{(l)}$ 
10    else
11       $r_j^{(l+1)} := r_j^{(l)} + m$ 
12       $\text{changed} := \text{true}$ 
13    end
14  end
15   $l := l + 1$ 
16 while  $\text{changed}$ 

```

To analyze the complexity of the algorithm, we assume for ease of notation that

$$r_1^{(0)} = \dots = r_{d-1}^{(0)} = r^{(0)} \quad \text{and} \quad \hat{r}_1 = \dots = \hat{r}_{d-1} = \hat{r} = r^{(0)} + km - q$$

for some integer k . The algorithm carries out approximately k iterations, with the l^{th} one (for $l \in \{1, \dots, k\}$) costing $\mathcal{O}(dnr^2(r^{(0)} + lm))$ flops because of Algorithm 11.

This results in an overall complexity equal to

$$\sum_{l=1}^k \mathcal{O}(dnr^2(r^{(0)} + lm)) = \mathcal{O}(dnr^2(r^{(0)} + km)k) = \mathcal{O}(dnr^2(\hat{r} + q)k)$$

We may expect to be able to optimize the algorithm in a recursive fashion by exploiting the fact that the tensors \mathcal{A}^l have ranks which increase as the algorithm proceeds. Indeed, we can reduce the total cost of the sampling step (see lines 2 and 6 in Algorithm 11) by recycling the previously used random vectors. However, it is impossible to speed up the computation of the compressed cores (line 10), since the orthogonal matrices W_j' (for $j \in \{2, \dots, d-1\}$) cannot be updated in an incremental way.

4.4 Numerical examples

4.4.1 Truncation of a Scholes-like tensor

We consider $d > 0$, a domain $\Omega \subset \mathbb{R}^d$, and a symmetric matrix $\Sigma \in \mathbb{R}^{d \times d}$ with entries $\Sigma(\alpha, \beta) = \sigma_{\alpha\beta}$ for $(\alpha, \beta) \in \{1, \dots, d\}^2$. We define a second order differential operator over $C^2(\Omega)$ as (see [PP07])

$$\mathcal{L} = \sum_{\alpha=1}^d \sum_{\beta=\alpha+1}^d \sigma_{\alpha\beta} \frac{\partial^2}{\partial x_\alpha \partial x_\beta}$$

which can be discretized, e.g. using finite differences, with the linear operator

$$A = \sum_{\alpha=1}^d \sum_{\beta=\alpha+1}^d \sigma_{\alpha\beta} W_\alpha W_\beta$$

In particular, each W_α (for $\alpha \in \{1, \dots, d\}$) can be interpreted as the discrete gradient operator with respect to the α^{th} component.

If Ω is discretized using a d -dimensional uniform Cartesian grid with m points in each dimension, the operators W_α and A can be interpreted as $m^d \times m^d$ matrices:

$$W_\alpha = I_m \otimes \dots \otimes I_m \otimes \underbrace{B}_{\alpha^{\text{th}} \text{ position}} \otimes I_m \otimes \dots \otimes I_m$$

and

$$A = \sum_{\alpha=1}^{d-1} \sum_{\beta=\alpha+1}^d \sigma_{\alpha\beta} \left(I_m \otimes \dots \otimes I_m \otimes \underbrace{B}_{\alpha^{\text{th}} \text{ position}} \otimes I_m \otimes \dots \right. \\ \left. \dots \otimes I_m \otimes \underbrace{B}_{\beta^{\text{th}} \text{ position}} \otimes I_m \otimes \dots \otimes I_m \right) \quad (4.25)$$

with B being a $m \times m$ discretization of the 1-dimensional gradient operator.

If each $m \times m$ matrix appearing in the Kronecker product is reshaped into a vector of size m^2 , using the change of indices

$$i^*(i, j) = i + m(j-1) \quad \text{for } (i, j) \in \{1, \dots, m\}^2,$$

then we can interpret \mathbf{A} as a $m^2 \times \dots \times m^2$ tensor \mathcal{A} of order d , whose entries are

$$\mathcal{A}(i^*(i_1, j_1), \dots, i^*(i_d, j_d)) = \mathbf{A}(m^{d-1}i_1 + \dots + mi_{d-1} + i_d, m^{d-1}j_1 + \dots + mj_{d-1} + j_d)$$

for $(i_l, j_l) \in \{1, \dots, m\}^2$, for $l \in \{1, \dots, d\}$.

In this representation, (4.25) can be viewed as the canonical decomposition of \mathcal{A} with canonical rank $d(d-1)/2$: indeed,

$$\mathcal{A}(i^*(i_1, j_1), \dots, i^*(i_d, j_d)) = \sum_{\alpha=1}^{d-1} \sum_{\beta=\alpha+1}^d \mathbf{U}_1((i_1, j_1), (\alpha, \beta)) \dots \mathbf{U}_d((i_d, j_d), (\alpha, \beta))$$

where

$$\mathbf{U}_l((i_l, j_l), (\alpha, \beta)) = \begin{cases} \mathbf{B}(i_l, j_l) & \text{if } l \in \{\alpha, \beta\} \\ \mathbf{I}_m(i_l, j_l) & \text{if } l \notin \{\alpha, \beta\} \end{cases} \quad \text{for } l \in \{1, \dots, d\}$$

From here, we can derive a representation of \mathcal{A} in the TT format with TT-ranks $(1, d(d-1)/2, \dots, d(d-1)/2, 1)$ by applying the strategy described in Section 4.1.

However, it can be proven (see [Ose11]) that there exists a TT representation of \mathcal{A} whose TT-ranks $(1, r_1^*, \dots, r_{d-1}^*, 1)$ are given by

$$r_j^* = 2 + \min\{j, d-j\} \quad \text{for } j \in \{1, \dots, d-1\}$$

For this reason, we wish to compress the tensor to TT-ranks $(1, r_1^*, \dots, r_{d-1}^*, 1)$ by applying Algorithms 7 and 11 (respectively deterministic and randomized truncation) for $d \in \{5, 10, \dots, 40\}$. We choose $m = 10$ and \mathbf{B} as the first-order finite differences approximation of the 1-dimensional gradient operator on a uniform grid (we ignore the leading constant due to the grid size, since it does not affect our analysis). The elements of Σ are randomly drawn from the uniform distribution.

The comparison is carried out in the computational environment described in Section 3.3.1, yielding the results shown in Figure 4.5. In particular, in this and the following examples, we employ the TTeMPS toolbox (<http://anchp.epfl.ch/TTeMPS>) for the storage and for most of the basic operations (e.g. sums, norms, orthogonalizations) in the TT format.

For all simulations, both algorithms yield a relative error smaller than 10^{-13} , in agreement with the theoretical bound on the TT-ranks. As d increases, the error tends to increase: for both algorithms this is mainly due to the round-off errors introduced by computing QR decompositions (and SVDs in the deterministic approach).

In particular, the error appears higher (in the majority of the simulations) for Algorithm 3: most likely, this is due to the fact that the deterministic algorithm has to compute QR decompositions of bigger matrices. Indeed, the first sweep of Algorithm 3 deals with matrices of size $(nd(d-1)/2) \times (d(d-1)/2)$, whereas the largest matrix whose QR decomposition is computed in Algorithm 7 has approximately size $(n(d/2 + p + 2)) \times (d/2 + p + 2)$.

In terms of computing time, the deterministic algorithm runs faster for small tensors ($d \in \{5, 10\}$), while the opposite holds true for bigger tensors. For instance, for $d = 40$, Algorithm 11 is more than four times faster than Algorithm 7, with a running time of 70.0 seconds versus 304 seconds.

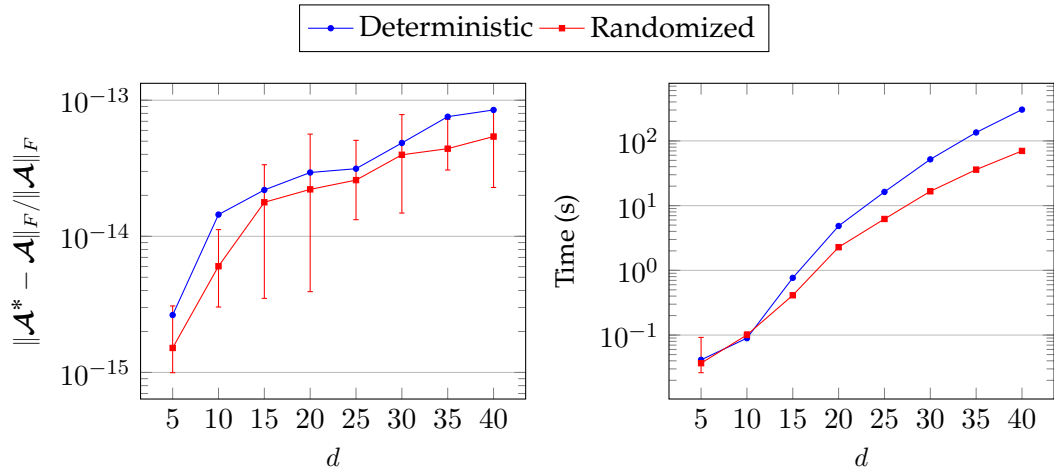


FIGURE 4.5: Relative error (left) and computing time (right) for Algorithms 7 and 11. For the randomized algorithm, we use an oversampling $p = 2$. The error bars show the range of values obtained over $N = 20$ simulations.

4.4.2 Truncation of the discretization of a rational function

We consider the positive integers d and n , as well as $\gamma \in \mathbb{R}^+$, and we discretize uniformly the interval $[\frac{1}{10}, \frac{n}{10}]$ with the points $x_i = \frac{i}{10}$, for $i \in \{1, \dots, n\}$. We define the d -dimensional tensor \mathcal{A} of size $n \times \dots \times n$ as the natural generalization of the matrix A defined in Section 3.3.2:

$$\mathcal{A}(i_1, \dots, i_d) = f_{d,\gamma}(x_{i_1}, \dots, x_{i_d}) = (x_{i_1} + \dots + x_{i_d})^{-\gamma} \quad \text{for } (i_1, \dots, i_d) \in \{1, \dots, n\}^d$$

Due to the smoothness of $f_{d,\gamma}$ in $(0, \frac{n+1}{10})^d$, each unfolding of \mathcal{A} has singular values which decay (at least) exponentially (see e.g. [GL17]). Indeed, if the tensor is compressed down to its numerical rank (i.e. it is rounded to machine precision), the ranks are quite small: for $n = 50$, $\gamma = 1$ and $d \in \{3, 5, 10\}$ the TT-ranks are uniformly bounded by 10.

Given the values of n , γ and d specified above, the tensor \mathcal{A} is generated by applying an exponential sum approximation (see e.g. [BH05] and [McL16]) with 25 terms. With this choice, we obtain a tensor with uniform TT-ranks $(1, 25, \dots, 25, 1)$, which approximates the function $f_{d,\gamma}$ with (at least) a uniform precision $2.3 \cdot 10^{-10}$ over the n^d sample points $\{\frac{1}{10}, \dots, \frac{n}{10}\}^d$.

Now we compress the tensor to uniform TT-ranks $(1, 4, \dots, 4, 1)$ using Algorithm 11, while keeping the result of Algorithm 7 as a benchmark. In particular, we are interested in observing how the results depend on the oversampling parameter p .

The behavior of the relative error is shown in Figure 4.6. As in the matrix case (see Section 3.3.2), small values of p (e.g. $p = 3$) are sufficient to achieve the optimal accuracy. This agrees with the fast decay rate of the singular values of all the unfoldings of \mathcal{A} .

In particular, the error appears to converge to the optimal value with (at least) exponential rate in p , even though the theoretical analysis of the randomized approach (see Section 3.1.1) cannot be easily generalized to “rank-1” random vectors.

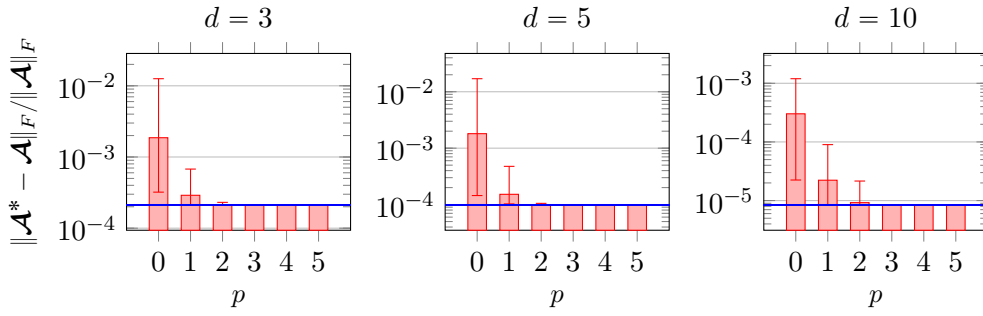


FIGURE 4.6: Relative error versus oversampling for Algorithm 11. The bars connect the minimum and the maximum value of the relative error obtained over $N = 20$ simulations. In blue the error obtained with the deterministic algorithm.

4.4.3 Rounding of a random tensor

Given the positive integers n, r and d (with $n \geq r$), we can build a d -dimensional tensor \mathcal{A} of size $n \times \dots \times n$ with TT-ranks $(1, r, \dots, r, 1)$ in the following way. We consider d randomly generated $n \times r$ matrices with orthogonal columns $\{U_j\}_{j=1}^d$ and we use them to prescribe the canonical decomposition of \mathcal{A} :

$$\mathcal{A}(i_1, \dots, i_d) = \sum_{\alpha=1}^r U_1(i_1, \alpha) \dots U_d(i_d, \alpha)$$

It is trivial to verify that the TT decomposition of \mathcal{A} obtained by using the method described in Section 4.1 is both left- and right-orthogonal.

With a small modification, we can prescribe an arbitrary decay for the singular values of the unfoldings. Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ and consider

$$\mathcal{B}(i_1, \dots, i_d) = \sum_{\alpha=1}^r \sigma_\alpha U_1(i_1, \alpha) \dots U_d(i_d, \alpha) \quad (4.26)$$

Proceeding again as in Section 4.1, we may compute a right-orthogonal TT decomposition of \mathcal{B} with

$$\mathcal{B}_1(1, i_1, \alpha_1) = \sigma_{\alpha_1} U_1(i_1, \alpha_1) \quad \text{for } (i_1, \alpha_1) \in I_1 \times \{1, \dots, r\}$$

For this decomposition, it is easy to verify that all unfoldings have (non-zero) singular values $(\sigma_1, \dots, \sigma_r)$.

Now we fix $n = 50, r = 50$ and $d = 20$, and we prescribe an exponential decay $\sigma_j = e^{1-j}$ for $j \in \{1, \dots, r\}$. Using Algorithms 8 and 12, we wish to compress \mathcal{B} down to precision $\varepsilon \in \{10^{-2}, 10^{-4}, 10^{-6}, 10^{-8}\}$. In particular, we choose $m = 3$ and $q = 2$ in the randomized approach.

By generalizing the Eckart-Young theorem (see e.g. [EY36]), we can identify the solution of the rounding problem whose TT-ranks are uniformly the smallest among all solutions: if we want the relative error (4.10) to be smaller than ε , it is sufficient to truncate the sum appearing in (4.26) at the term r^* , with r^* being the smallest positive integer such that

$$\sum_{j=r^*+1}^r \sigma_j^2 \leq \varepsilon^2 \sum_{j=1}^r \sigma_j^2$$

For the choices of $(\sigma_j)_{j=1}^r$ and ε specified above, the corresponding values of r^* are $\{6, 11, 15, 20\}$ respectively.

For the randomized algorithm we consider two cases: first, we assume that no initial guess is given, or equivalently that $r_1^{(0)} = \dots = r_{d-1}^{(0)} = 0$; secondly, we assume that we are given an initial guess with $r_1^{(0)} = \dots = r_{d-1}^{(0)} = \lceil r^*/2 \rceil$.

In Figure 4.7 we show the results of the numerical experiment, which is carried out in the environment described in Section 3.3.1. From the left plot, we can observe that the required precision is almost always satisfied, and that the relative error is within a small factor of ε even in case of failure. From the right plot we can see that Algorithm 12 is always faster than Algorithm 8. We can also observe that the gap between the speed of the two approaches is particularly pronounced for larger values of ε : this is due to the fact that the cost of the randomized approach depends only weakly on the prescribed precision.

In Figure 4.8 we report the largest TT-ranks of the compressed tensors at the end of each algorithm. For the randomized approach we also plot the values before the last recompression, i.e. the final TT-ranks of \mathcal{C}' in Algorithm 12. In all cases, we can verify that only a modest oversampling (measured as the variation between the TT-ranks before and after recompression) is needed for the completion of Algorithm 12.

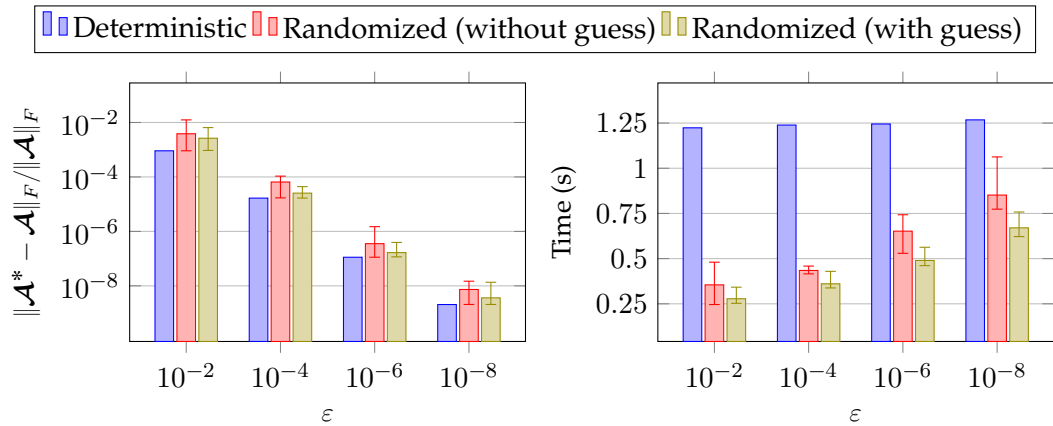


FIGURE 4.7: Relative error (left) and computing time (right) for Algorithms 8 and 12. The bars connect the minimum and the maximum values obtained over $N = 20$ simulations.

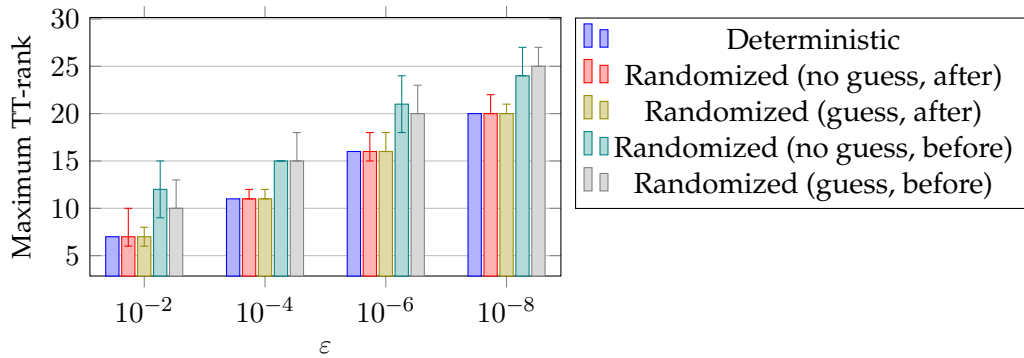


FIGURE 4.8: Maximum TT-rank of the compressed tensor for Algorithms 8 and 12. For the randomized algorithm the plot shows the modal values of the maximal TT-ranks over $N = 20$ simulations, both before and after recompression. Also, the bars show the range of values obtained.

Chapter 5

Compression of Hadamard products

Hadamard products are a fundamental operation in matrix and tensor computations, whose importance is mostly due to the fact that they represent the discrete version of point-wise products of functions. They can be carried out in a very trivial way on matrices and tensors which are stored as full arrays.

However, many issues arise if the two factors are in a structured format, e.g. the low-rank dyadic decomposition (1.1) for matrices and the TT format (4.2) for tensors. In these cases, Hadamard products increase the sizes of the factors/cores: as we show in the upcoming sections, all the ranks grow in a quadratic fashion. This leads to a substantial increase in the computational cost of most manipulations of the matrix/tensor, and also in the memory required for its storage.

For this reason, it is important to recompress properly the result of the Hadamard product. Still, methods which rely, as an intermediate step, on the computation of the full product are quite expensive, since they are affected by the issues described above. In this chapter, we devise randomized algorithms which can compress Hadamard products without computing the full matrix/tensor, exploiting the particular structure that Hadamard products induce on the dyadic decomposition and on the TT format.

5.1 Compression of Hadamard products of matrices

Consider two $m \times n$ matrices \mathbf{A} and \mathbf{B} , whose low-rank approximations, with ranks r_A and r_B respectively, are known:

$$\mathbf{A} = \mathbf{W}_A \mathbf{Z}_A^\top \quad \text{and} \quad \mathbf{B} = \mathbf{W}_B \mathbf{Z}_B^\top$$

In particular, we assume that $r_A \geq r_B$ without loss of generality.

A rank- $(r_A r_B)$ dyadic decomposition of $\mathbf{C} = \mathbf{A} * \mathbf{B}$ can be found quite easily:

$$\mathbf{C} = (\mathbf{W}_A \odot^\top \mathbf{W}_B) (\mathbf{Z}_A \odot^\top \mathbf{Z}_B)^\top = (\mathbf{W}_A \odot^\top \mathbf{W}_B) (\mathbf{Z}_A^\top \odot \mathbf{Z}_B^\top) \quad (5.1)$$

Indeed, for any $(i, j) \in \{1, \dots, m\} \times \{1, \dots, n\}$, it holds

$$C(i, j) = A(i, j)B(i, j) = \sum_{l=1}^{r_A} \sum_{k=1}^{r_B} W_A(i, l) Z_A(j, l) W_B(i, k) Z_B(j, k)$$

$$\begin{aligned}
&= \sum_{l=1}^{r_A} \sum_{k=1}^{r_B} (\mathbf{W}_A \odot^\top \mathbf{W}_B) (i, k + (l-1)r_B) (\mathbf{Z}_A \odot^\top \mathbf{Z}_B) (j, k + (l-1)r_B) \\
&= (\mathbf{W}_A \odot^\top \mathbf{W}_B) (i, :) (\mathbf{Z}_A^\top \odot \mathbf{Z}_B^\top) (:, j)
\end{aligned}$$

In order to recompress \mathbf{C} to a lower rank, we may follow a deterministic approach, as follows:

- Compute the QR decomposition of the right factor: $\mathbf{Z}_A \odot^\top \mathbf{Z}_B = \mathbf{Q}\mathbf{R}$.
- Update the left factor: $\mathbf{W}' = (\mathbf{W}_A \odot^\top \mathbf{W}_B) \mathbf{R}^\top$.
- Obtain a low-rank approximation of \mathbf{W}' with the desired rank or precision (e.g. through a RRQR decomposition):

$$\mathbf{W}' \approx \tilde{\mathbf{W}} \tilde{\mathbf{Z}}^\top$$

- Define the new low-rank factors of \mathbf{C} as $\mathbf{W}^* = \tilde{\mathbf{W}}$ and $\mathbf{Z}^* = \mathbf{Q}\tilde{\mathbf{Z}}$

Instead, a randomized approach can obtain a low-rank approximation of \mathbf{C} without computing the full factors $\mathbf{W}_A \odot^\top \mathbf{W}_B$ and $\mathbf{Z}_A \odot^\top \mathbf{Z}_B$. To this aim, we modify the algorithms described in Section 3.1 to exploit the particular structure of \mathbf{C} : indeed, we can compute matrix-vector products involving \mathbf{C} in $\mathcal{O}(mr_A r_B)$ flops, using only $\mathcal{O}(mr_B + nr_A)$ extra memory, as follows.

For $i \in \{1, \dots, m\}$, it holds that

$$\begin{aligned}
(\mathbf{C}\mathbf{x})(i) &= \sum_{j=1}^n \mathbf{A}(i, j) \mathbf{B}(i, j) \mathbf{x}(j) = \\
&= \sum_{j,k=1}^n \mathbf{A}(i, j) \mathbf{B}(i, k) \mathbf{x}(j) \mathbf{I}_n(j, k) = (\mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}^\top) (i, i)
\end{aligned}$$

which implies that

$$\mathbf{C}\mathbf{x} = \operatorname{diag}(\mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}^\top) = \operatorname{diag}(\mathbf{W}_A \mathbf{Z}_A^\top \operatorname{diag}(\mathbf{x}) \mathbf{Z}_B \mathbf{W}_B^\top) \quad (5.2)$$

Thus, we can proceed with the computation from inside out:

- $\mathbf{Y}_1 = \operatorname{diag}(\mathbf{x}) \mathbf{Z}_B$ can be computed in just $\mathcal{O}(nr_B)$ flops due to the diagonal structure of the first factor. Indeed, it suffices to evaluate each row as

$$\mathbf{Y}_1(i, :) = \mathbf{x}(i) \mathbf{Z}_B(i, :) \quad \text{for } i \in \{1, \dots, n\}$$

- $\mathbf{Y}_2 = \mathbf{W}_A (\mathbf{Z}_A^\top \mathbf{Y}_1)$ can be computed in $\mathcal{O}(mr_A r_B)$ flops.
- Since we just need the diagonal terms of $\mathbf{Y}_2 \mathbf{W}_B^\top$, we do not compute the full matrix. Instead, we can just obtain each component as

$$(\mathbf{C}\mathbf{x})(i) = \sum_{j=1}^{r_B} \mathbf{Y}_2(i, j) \mathbf{W}_B(i, j) \quad \text{for } i \in \{1, \dots, m\}$$

using $\mathcal{O}(mr_B)$ flops overall.

Similarly, the computation of matrix-vector products involving C^\top can be carried out as:

$$C^\top \mathbf{y} = \text{diag}(\mathbf{A}^\top \text{diag}(\mathbf{y}) \mathbf{B}) = \text{diag}(\mathbf{Z}_A \mathbf{W}_A^\top \text{diag}(\mathbf{y}) \mathbf{W}_B \mathbf{Z}_B^\top) \quad (5.3)$$

Again, this requires only $\mathcal{O}(mr_A r_B)$ flops, and $\mathcal{O}(mr_B + nr_A)$ extra memory.

Using (5.2) and (5.3), we can specialize Algorithms 1 and 3 to tackle the case of Hadamard products. We report the extension of Algorithm 3 in Algorithm 13. The modified version of Algorithm 1 can be obtained in the same spirit.

Algorithm 13: Randomized compression with fixed (Frobenius) precision

Data: matrices $\mathbf{W}_A \in \mathbb{R}^{m \times r_A}$, $\mathbf{Z}_A \in \mathbb{R}^{n \times r_A}$, $\mathbf{W}_B \in \mathbb{R}^{m \times r_B}$, and $\mathbf{Z}_B \in \mathbb{R}^{n \times r_B}$,
 $\varepsilon > 0$, integer $q > 0$

Result: $\mathbf{W} \in \mathbb{R}^{m \times r}$ and $\mathbf{Z} \in \mathbb{R}^{n \times r}$, with $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_m$ and such that

$$\|(\mathbf{W}_A \mathbf{Z}_A^\top) * (\mathbf{W}_B \mathbf{Z}_B^\top) - \mathbf{W} \mathbf{Z}^\top\|_F \leq \varepsilon \|(\mathbf{W}_A \mathbf{Z}_A^\top) * (\mathbf{W}_B \mathbf{Z}_B^\top)\|_F$$

holds with high probability

```

1 Generate Gaussian matrix  $\Omega \in \mathbb{R}^{n \times q}$ 
2  $\mathbf{Y} := ((\mathbf{W}_A \mathbf{Z}_A^\top) * (\mathbf{W}_B \mathbf{Z}_B^\top)) \Omega$  % column-wise using (5.2)
3 Compute QR decomposition:  $\mathbf{Y} =: \mathbf{W} \mathbf{R}$ 
4  $\mathbf{Z} := ((\mathbf{Z}_A \mathbf{W}_A^\top) * (\mathbf{Z}_B \mathbf{W}_B^\top)) \mathbf{W}$  % column-wise using (5.3)
5  $r := 0$ 
6 while  $\|\mathbf{Z}\|_F^2 - \|\mathbf{Z}(:, 1:r)\|_F^2 > \varepsilon^2 \|\mathbf{Z}\|_F^2$  do
7    $r := r + 1$ 
8   Generate Gaussian vector  $\omega \in \mathbb{R}^n$ 
9    $\mathbf{y} := ((\mathbf{W}_A \mathbf{Z}_A^\top) * (\mathbf{W}_B \mathbf{Z}_B^\top)) \omega$  % using (5.2)
10   $\mathbf{w} := (\mathbf{I}_m - \mathbf{W} \mathbf{W}^\top) \mathbf{y}$ 
11   $\mathbf{w} := \mathbf{w} / \|\mathbf{w}\|$ 
12   $\mathbf{W} := [\mathbf{W} \mid \mathbf{w}]$ 
13   $\mathbf{z} := ((\mathbf{Z}_A \mathbf{W}_A^\top) * (\mathbf{Z}_B \mathbf{W}_B^\top)) \mathbf{w}$  % using (5.3)
14   $\mathbf{Z} := [\mathbf{Z} \mid \mathbf{z}]$ 
15 end

```

The complexity of Algorithm 13 is $\mathcal{O}(mr_A r_B (r + q))$, with the dominant steps being the matrix-vector products carried out by exploiting (5.2) and (5.3).

Moreover, the total memory required throughout the execution of the algorithm is just $\mathcal{O}(mr_B + nr_A + m(r + q))$, due to (5.2) and (5.3), and to the storage of the final factors. Under the reasonable (for ε small enough) assumption that $r \geq r_A$, it can be simplified as $\mathcal{O}(m(r + q))$.

Similarly, the extension of the fixed-rank case requires $\mathcal{O}(mr_A r_B (r + p))$ flops, with p being the oversampling parameter, and $\mathcal{O}(m(r + p))$ memory.

In comparison, the deterministic approach described above costs approximately $\mathcal{O}(mr_A^2 r_B^2)$ flops, and needs $\mathcal{O}(mr_A r_B)$ memory.

5.2 Compression of Hadamard products of tensors in the TT format

Consider two tensors of size $n_1 \times \dots \times n_d$ in the TT format:

$$\mathcal{A}(i_1, \dots, i_d) = \mathcal{A}_1(i_1) \dots \mathcal{A}_d(i_d) \quad \text{and} \quad \mathcal{B}(i_1, \dots, i_d) = \mathcal{B}_1(i_1) \dots \mathcal{B}_d(i_d)$$

for $(i_1, \dots, i_d) \in I_1 \times \dots \times I_d$, with TT-ranks $(1, r_1^A, \dots, r_{d-1}^A, 1)$ and $(1, r_1^B, \dots, r_{d-1}^B, 1)$ respectively.

Using the same reasoning applied in the previous section, we can easily verify (see e.g. [Ose11]) that the tensor $\mathcal{C} = \mathcal{A} * \mathcal{B}$ admits the following TT representation with TT-ranks $(1, r_1^A r_1^B, \dots, r_{d-1}^A r_{d-1}^B, 1)$:

$$\mathcal{C}(i_1, \dots, i_d) = \underbrace{(\mathcal{A}_1(i_1) \otimes \mathcal{B}_1(i_1))}_{\mathbf{c}_1(i_1)} \dots \underbrace{(\mathcal{A}_d(i_d) \otimes \mathcal{B}_d(i_d))}_{\mathbf{c}_d(i_d)} \quad (5.4)$$

for $(i_1, \dots, i_d) \in I_1 \times \dots \times I_d$.

As before, we wish to recompress \mathcal{C} down to some prescribed rank or precision. One way to do this is to explicitly compute each of the $(r_{j-1}^A r_{j-1}^B) \times n_j \times (r_j^A r_j^B)$ cores appearing in (5.4) (for $j = 1, \dots, d$), and then apply Algorithm 7 or 8, depending on the kind of compression we wish to achieve.

In particular, it is important to observe the following property: even if the TT decompositions of both \mathcal{A} and \mathcal{B} are (either left- or right-) orthogonal, in general the TT decomposition (5.4) is not.

Indeed, assume that both \mathcal{A} and \mathcal{B} have left-orthogonal TT decompositions. For any $j \in \{1, \dots, d\}$ it holds

$$\begin{aligned} (\mathbf{c}_j^{<2>})^\top \mathbf{c}_j^{<2>} &= \sum_{i_j=1}^{n_j} \mathbf{c}_j(i_j)^\top \mathbf{c}_j(i_j) \\ &= \sum_{i_j=1}^{n_j} (\mathcal{A}_j(i_j) \otimes \mathcal{B}_j(i_j))^\top (\mathcal{A}_j(i_j) \otimes \mathcal{B}_j(i_j)) \\ &= \sum_{i_j=1}^{n_j} (\mathcal{A}_j(i_j)^\top \mathcal{A}_j(i_j)) \otimes (\mathcal{B}_j(i_j)^\top \mathcal{B}_j(i_j)) \end{aligned}$$

which, in general, is different from

$$\begin{aligned} \sum_{i_j=1}^{n_j} \sum_{i'_j=1}^{n_j} (\mathcal{A}_j(i_j)^\top \mathcal{A}_j(i_j)) \otimes (\mathcal{B}_j(i'_j)^\top \mathcal{B}_j(i'_j)) &= \\ &= \left(\sum_{i_j=1}^{n_j} \mathcal{A}_j(i_j)^\top \mathcal{A}_j(i_j) \right) \otimes \left(\sum_{i'_j=1}^{n_j} \mathcal{B}_j(i'_j)^\top \mathcal{B}_j(i'_j) \right) = \\ &= \left((\mathcal{A}_j^{<2>})^\top \mathcal{A}_j^{<2>} \right) \otimes \left((\mathcal{B}_j^{<2>})^\top \mathcal{B}_j^{<2>} \right) = \mathbf{I}_{r_j^A} \otimes \mathbf{I}_{r_j^B} = \mathbf{I}_{r_j^A r_j^B} \end{aligned}$$

As such, it is impossible to avoid the preliminary orthogonalization step in Algorithms 7 and 8. This leads us to conclude that the deterministic approach is quite

expensive: assuming that

$$r_1^A = \dots = r_{d-1}^A = r_1^B = \dots = r_{d-1}^B = r \quad \text{and} \quad n_1 = \dots = n_d = n \quad (5.5)$$

the procedure costs $\mathcal{O}(dnr^6)$ flops (due to the orthogonalization step), and needs $\mathcal{O}(dnr^4)$ memory (to store \mathcal{C} before compressing it), independently of the final TT-ranks.

As in the matrix case, it is possible to improve both complexity and memory requirements using a randomized approach.

5.2.1 Sampling the range of the unfoldings of Hadamard products

First, we consider the problem of finding a low-rank approximation for the j^{th} unfolding of \mathcal{C} , with $j \in \{1, \dots, d-1\}$. As in Section 4.3, we solve this task by sampling the range of the unfolding with some random vectors:

$$\mathbf{y}^{(\alpha)} = \mathcal{C}^{<j>} \boldsymbol{\omega}^{(\alpha)} \quad \text{with} \quad \boldsymbol{\omega}^{(\alpha)} \in \mathbb{R}^{(n_{j+1} \dots n_d)}, \quad \text{for} \quad \alpha \in \{1, 2, \dots\}$$

In particular, we have seen that, if we choose “rank-1” random vectors of the form

$$\boldsymbol{\omega}^{(\alpha)} = \boldsymbol{\omega}_d^{(\alpha)} \otimes \dots \otimes \boldsymbol{\omega}_{j+1}^{(\alpha)}, \quad \text{with} \quad \boldsymbol{\omega}_l^{(\alpha)} \in \mathbb{R}^{n_l} \quad \text{for} \quad l \in \{j+1, \dots, d\},$$

then we can use multidimensional contractions to compute efficiently the samples. In addition, in this case we can exploit the structure of \mathcal{C} to speed up the process even further, and avoid storing the uncompressed cores \mathcal{C}_l , for $l \in \{j+1, \dots, d\}$.

Indeed, let $l \in \{j+1, \dots, d\}$ and consider lines 6–8 in Algorithm 10: in the case at hand, we want to compute

$$\mathbf{y}_l^{(\alpha)} = \mathcal{C}_{\geq l}^\top \left(\boldsymbol{\omega}_d^{(\alpha)} \otimes \dots \otimes \boldsymbol{\omega}_l^{(\alpha)} \right) = \sum_{i_l=1}^{n_l} \left(\mathcal{C}_l(i_l) \mathbf{y}_{l+1}^{(\alpha)} \right) \boldsymbol{\omega}_l^{(\alpha)}(i_l)$$

with $\mathbf{y}_{l+1}^{(\alpha)} \in \mathbb{R}^{(r_l^A r_l^B)}$ being either a previously computed vector, or 1 if $l = d$. Due to (5.4), this can be expressed as

$$\mathbf{y}_l^{(\alpha)} = \sum_{i_l=1}^{n_l} \left((\mathcal{A}_l(i_l) \otimes \mathcal{B}_l(i_l)) \mathbf{y}_{l+1}^{(\alpha)} \right) \boldsymbol{\omega}_l^{(\alpha)}(i_l) \quad (5.6)$$

Now, consider the matrices $\mathbf{Y}_l^{(\alpha)} \in \mathbb{R}^{r_{l-1}^B \times r_{l-1}^A}$ and $\mathbf{Y}_{l+1}^{(\alpha)} \in \mathbb{R}^{r_l^B \times r_l^A}$, obtained by reshaping the vectors $\mathbf{y}_l^{(\alpha)}$ and $\mathbf{y}_{l+1}^{(\alpha)}$ respectively, so that

$$\left(\mathbf{Y}_l^{(\alpha)} \right)^{<2>} = \mathbf{y}_l^{(\alpha)} \quad \text{and} \quad \left(\mathbf{Y}_{l+1}^{(\alpha)} \right)^{<2>} = \mathbf{y}_{l+1}^{(\alpha)}$$

We can exploit the fact that

$$(\mathbf{A} \otimes \mathbf{B}) (\mathbf{C})^{<2>} = (\mathbf{B} \mathbf{C} \mathbf{A}^\top)^{<2>} \quad (5.7)$$

which holds for any matrices \mathbf{A} , \mathbf{B} and \mathbf{C} of suitable size (see e.g. [Sch04]), and rewrite (5.6) as

$$\mathbf{Y}_l^{(\alpha)} = \sum_{i_l=1}^{n_l} \left(\mathbf{B}_l(i_l) \mathbf{Y}_{l+1}^{(\alpha)} \mathbf{A}_l(i_l)^\top \right) \boldsymbol{\omega}_l^{(\alpha)}(i_l) \quad (5.8)$$

In this way we can compute the new sample in $\mathcal{O}(n_l r_{l-1}^A r_l^B (r_l^A + r_{l-1}^B))$ flops at most.

Hence, Algorithm 10 can be efficiently applied to the case of Hadamard products. If we employ assumption (5.5), the total cost of generating a single set of samples (i.e. running the algorithm with $k = 1$) is $\mathcal{O}(dnr^3)$ flops, and the memory required (excluding the storage of the tensors \mathbf{A} and \mathbf{B}) is just $\mathcal{O}(dr^2)$.

5.2.2 Low-rank approximation of tensors in the TT format

Let us consider the case of the truncation of the Hadamard product of two tensors to some prescribed TT-ranks $(1, r_1^*, \dots, r_{d-1}^*, 1)$ (the rounding case will follow as a generalization of Algorithm 12).

First, we run Algorithm 10 as discussed in the previous section. In particular, we set $k = \max_{j=1, \dots, d-1} r_j^* + p$, with $p > 0$ being the oversampling parameter. In this way, we obtain k samples from the row space of each right interface matrix $\mathbf{C}_{\geq j}$, for $j \in \{2, \dots, d\}$. As before, we reshape the samples as matrices:

$$\mathbf{Y}_j^{(l)} \in \mathbb{R}^{r_{j-1}^B \times r_{j-1}^A}, \quad \text{with } \left(\mathbf{Y}_j^{(l)} \right)^{\langle 2 \rangle} = \mathbf{C}_{\geq j}^\top \left(\boldsymbol{\omega}_d^{(l)} \otimes \dots \otimes \boldsymbol{\omega}_j^{(l)} \right)$$

for $l \in \{1, \dots, k\}$ and $j \in \{2, \dots, d\}$.

Now we can apply a procedure similar to Algorithm 11, and iteratively compute each core of the compressed tensor, as follows. The first unfolding of \mathbf{C} is of the form

$$\mathbf{C}^{\langle 1 \rangle} = \mathbf{C}_{\leq 1} \mathbf{C}_{\geq 2}^\top = (\mathbf{C}_1)^{\langle 2 \rangle} \mathbf{C}_{\geq 2}^\top$$

and its range can be sampled as

$$\mathbf{C}^{\langle 1 \rangle} \boldsymbol{\omega}^{(l)} = \mathbf{C}^{\langle 1 \rangle} \left(\boldsymbol{\omega}_d^{(l)} \otimes \dots \otimes \boldsymbol{\omega}_2^{(l)} \right) = (\mathbf{C}_1)^{\langle 2 \rangle} \left(\mathbf{Y}_2^{(l)} \right)^{\langle 2 \rangle}$$

for $l \in \{1, \dots, r_1^* + p\}$.

Due to (5.4), we have the following representation of each sample:

$$\mathbf{C}^{\langle 1 \rangle} \boldsymbol{\omega}^{(l)} = \begin{bmatrix} \mathbf{A}_1(1) \otimes \mathbf{B}_1(1) \\ \vdots \\ \mathbf{A}_1(n_1) \otimes \mathbf{B}_1(n_1) \end{bmatrix} \left(\mathbf{Y}_2^{(l)} \right)^{\langle 2 \rangle} = \begin{bmatrix} (\mathbf{A}_1(1) \otimes \mathbf{B}_1(1)) \left(\mathbf{Y}_2^{(l)} \right)^{\langle 2 \rangle} \\ \vdots \\ (\mathbf{A}_1(n_1) \otimes \mathbf{B}_1(n_1)) \left(\mathbf{Y}_2^{(l)} \right)^{\langle 2 \rangle} \end{bmatrix}$$

Using (5.7), this can be expressed as

$$\mathbf{C}^{\langle 1 \rangle} \boldsymbol{\omega}^{(l)} = \begin{bmatrix} \left(\mathbf{B}_1(1) \mathbf{Y}_2^{(l)} \mathbf{A}_1(1)^\top \right)^{\langle 2 \rangle} \\ \vdots \\ \left(\mathbf{B}_1(n_1) \mathbf{Y}_2^{(l)} \mathbf{A}_1(n_1)^\top \right)^{\langle 2 \rangle} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1(1) \mathbf{Y}_2^{(l)} \mathbf{A}_1(1)^\top \\ \vdots \\ \mathbf{B}_1(n_1) \mathbf{Y}_2^{(l)} \mathbf{A}_1(n_1)^\top \end{bmatrix} \quad (5.9)$$

where the “second unfolding” operator has been omitted since it acts on a scalar.

We can merge all the samples into the $n_1 \times (r_1^* + p)$ matrix

$$\mathbf{C}^{<1>} \left[\boldsymbol{\omega}^{(1)} \mid \dots \mid \boldsymbol{\omega}^{(r_1^*+p)} \right]$$

and compute its QR decomposition. If we indicate with \mathbf{W}_1 the orthogonal factor of such a decomposition, we have the following rank- $(r_1^* + p)$ approximation for the first unfolding of \mathbf{C} :

$$\mathbf{C}^{<1>} \approx \mathbf{W}_1 \left(\mathbf{W}_1^\top (\mathbf{C}_1)^{<2>} \mathbf{C}_{\geq 2}^\top \right)$$

Now we can reshape \mathbf{W}_1 to obtain the first core of the compressed tensor:

$$\mathbf{C}'_1 \in \mathbb{R}^{1 \times n_1 \times (r_1^*+p)} \quad \text{with} \quad (\mathbf{C}'_1)^{<2>} = \mathbf{W}_1$$

Moreover, it is useful to define the $(r_1^A r_1^B) \times (r_1^* + p)$ “local interface matrix”

$$\mathbf{Z}_1 = (\mathbf{C}_1)^{<2>}^\top \mathbf{W}_1$$

which can be expressed as

$$\begin{aligned} \mathbf{Z}_1 &= [(\mathcal{A}_1(1) \otimes \mathcal{B}_1(1))^\top \mid \dots \mid (\mathcal{A}_1(n_1) \otimes \mathcal{B}_1(n_1))^\top] \mathbf{W}_1 \\ &= \sum_{i_1=1}^{n_1} (\mathcal{A}_1(i_1)^\top \otimes \mathcal{B}_1(i_1)^\top) \underbrace{\mathbf{W}_1(i_1, :)}_{\text{as a } 1 \times (r_1^*+p) \text{ matrix}} \end{aligned}$$

With these definitions, the second unfolding can be approximated as

$$\begin{aligned} \mathbf{C}^{<2>} &= \mathbf{C}_{\leq 2} \mathbf{C}_{\geq 3}^\top = \left(\mathbf{I}_{n_2} \otimes (\mathbf{C}_1)^{<2>} \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top \\ &\approx \left(\mathbf{I}_{n_2} \otimes (\mathbf{W}_1 \mathbf{Z}_1^\top) \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top \\ &= \left(\mathbf{I}_{n_2} \otimes \mathbf{W}_1 \right) \left(\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top \end{aligned}$$

Since \mathbf{W}_1 has orthogonal columns, $\mathbf{I}_{n_2} \otimes \mathbf{W}_1$ does as well, and it suffices to approximate the range of

$$\left(\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top$$

which can be sampled as

$$\left(\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top \boldsymbol{\omega}^{(l)} = \left(\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top \right) (\mathbf{C}_2)^{<2>} \left(\mathbf{Y}_3^{(l)} \right)^{<2>}$$

for $l \in \{1, \dots, r_2^* + p\}$.

Expression (5.9) can be extended to this case as well, yielding the block representation (of size $n_2(r_1^* + p)$)

$$\left(\mathbf{I}_{n_2} \otimes \mathbf{Z}_1^\top \right) (\mathbf{C}_2)^{<2>} \mathbf{C}_{\geq 3}^\top \boldsymbol{\omega}^{(l)} = \begin{bmatrix} \mathbf{Z}_1^\top \left(\mathcal{B}_2(1) \mathbf{Y}_3^{(l)} \mathcal{A}_2(1)^\top \right)^{<2>} \\ \vdots \\ \mathbf{Z}_1^\top \left(\mathcal{B}_2(n_2) \mathbf{Y}_3^{(l)} \mathcal{A}_2(n_2)^\top \right)^{<2>} \end{bmatrix} \quad (5.10)$$

for $l \in \{1, \dots, r_2^* + p\}$.

After repeating this operation for all the samples, we build a $(n_2(r_1^* + p)) \times (r_2^* + p)$ matrix, whose QR decomposition yields the orthogonal factor \mathbf{W}_2 . By reshaping this factor we obtain the second core of the compressed tensor

$$\mathbf{C}'_2 \in \mathbb{R}^{(r_1^* + p) \times n_2 \times (r_2^* + p)} \quad \text{with} \quad (\mathbf{C}'_2)^{\langle 2 \rangle} = \mathbf{W}_2$$

and we also compute the second “local interface matrix”

$$\mathbf{Z}_2 = (\mathbf{C}'_2)^{\langle 2 \rangle \top} (\mathbf{I}_{n_2} \otimes \mathbf{Z}_1) \mathbf{W}_2 \in \mathbb{R}^{(r_2^A r_2^B) \times (r_2^* + p)}$$

which, by generalizing (5.10), can be expressed as

$$\begin{aligned} \mathbf{Z}_2 &= [(\mathcal{A}_2(1) \otimes \mathcal{B}_1(1))^\top \mathbf{Z}_1 \mid \dots \mid (\mathcal{A}_2(n_2) \otimes \mathcal{B}_2(n_2))^\top \mathbf{Z}_1] \mathbf{W}_2 \\ &= \sum_{i_2=1}^{n_2} (\mathcal{A}_2(i_2)^\top \otimes \mathcal{B}_2(i_2)^\top) \mathbf{Z}_1 \mathbf{W}_2^{(i_2)} \end{aligned} \quad (5.11)$$

with $\mathbf{W}_2^{(i_2)} = \mathbf{W}_2((i_2 - 1)(r_1^* + p) + (1 : (r_1^* + p)), :)$ being a $(r_1^* + p) \times (r_2^* + p)$ block of \mathbf{W}_2 :

$$\mathbf{W}_2 = \begin{bmatrix} \mathbf{W}_2^{(1)} \\ \dots \\ \mathbf{W}_2^{(n_2)} \end{bmatrix}$$

Then we proceed iteratively until we compute the $(d - 1)$ th core \mathbf{C}'_{d-1} and the last “local interface matrix”

$$\mathbf{Z}_{d-1} = (\mathbf{C}'_{d-1})^{\langle 2 \rangle \top} (\mathbf{I}_{n_{d-1}} \otimes \mathbf{Z}_{d-2}) \mathbf{W}_{d-1} \in \mathbb{R}^{(r_{d-1}^A r_{d-1}^B) \times (r_{d-1}^* + p)}$$

which we use to compute the last core of the compressed tensor:

$$\mathbf{C}'_d = \mathbf{C}_d \times_1 \mathbf{Z}_{d-1}^\top \in \mathbb{R}^{(r_{d-1}^* + p) \times n_d \times 1}$$

Since the tensor \mathbf{C}' has TT-ranks $(1, r_1^* + p, \dots, r_{d-1}^* + p, 1)$, which are larger than the desired ones, we carry out a deterministic compression step by applying Algorithm 7. The complete procedure is summarized in Algorithm 14.

To analyze the complexity of the algorithm, it is necessary to identify the cost of lines 6 and 10. In particular, to simplify the notation, we assume that $r_1^* = \dots = r_{d-1}^* = r^*$.

We start from the former step, by considering the block decomposition introduced in (5.10). Under assumption (5.5), we can easily see that the computation of each of the n blocks requires $\mathcal{O}(r^3 + r^2(r^* + p))$ flops.

Similarly, the cost of line 10 can be found by referring to (5.11), and considering each term of the sum. First, we look at the product $(\mathcal{A}_j(i_j)^\top \otimes \mathcal{B}_j(i_j)^\top) \mathbf{Z}_{j-1}$. If we compute the product column-wise, we can apply (5.7) to carry out the whole multiplication in $\mathcal{O}(r^3(r^* + p))$ flops. Then we need to compute the product between the result and one of the blocks of \mathbf{W}_j , which can be done in $\mathcal{O}(r^2(r^* + p)^2)$ flops.

All the other operations appearing in Algorithm 14 are cheaper, or cost a similar number of flops. In particular, line 2 requires $\mathcal{O}(dnr^3(r^* + p))$ flops (see Section 5.2.1),

Algorithm 14: Randomized truncation of the Hadamard product of two tensors in the TT format

Data: tensors $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$ and $\mathcal{B}_1, \dots, \mathcal{B}_d$ respectively, desired TT-ranks $(r_0^* = 1, r_1^*, \dots, r_{d-1}^*, r_d^* = 1)$, integer $p \geq 0$

Result: $\mathcal{C}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{C}_1^*, \dots, \mathcal{C}_d^*$ with the specified TT-ranks, with $\mathcal{A} * \mathcal{B} \approx \mathcal{C}^*$

```

1  $k := \max\{r_1^*, \dots, r_{d-1}^*\} + p$ 
2  $\{\mathbf{Y}_2^{(1)}, \dots, \mathbf{Y}_2^{(k)}, \mathbf{Y}_3^{(1)}, \dots, \mathbf{Y}_d^{(k)}\} := \text{Algorithm 10}(\mathcal{C}, k)$  % using (5.8)
3  $\mathbf{Z}_0 := 1$ 
4 for  $j = 1, 2, \dots, d - 1$  do
5   for  $l = 1, 2, \dots, r_j^* + p$  do
6      $\mathbf{y}_j^{(l)} := (\mathbf{I}_{n_j} \otimes \mathbf{Z}_{j-1}^\top) (\mathcal{C}_j)^{\langle 2 \rangle} (\mathbf{Y}_{j-1}^{(l)})^{\langle 2 \rangle}$  % using (5.10)
7   end
8   Compute QR decomposition:  $[\mathbf{y}_j^{(1)} \mid \dots \mid \mathbf{y}_j^{(r_j^* + p)}] =: \mathbf{W}_j \mathbf{R}_j$ 
9    $\mathcal{C}'_j := \text{reshape}(\mathbf{W}_j, [r_{j-1}^* + p, n_j, r_j^* + p])$ 
10   $\mathbf{Z}_j := (\mathcal{C}_j)^{\langle 2 \rangle \top} (\mathbf{I}_{n_j} \otimes \mathbf{Z}_{j-1}) \mathbf{W}_j$  % using (5.11)
11 end
12  $\mathcal{C}'_d := \mathcal{C}_d \times_1 \mathbf{Z}_{d-1}^\top$ 
13  $\mathcal{C}^* := \text{Algorithm 7}(\mathcal{C}', (1, r_1^*, \dots, r_{d-1}^*, 1))$  % from right to left
```

and both the QR decompositions and the recompression step take $\mathcal{O}(dn(r^* + p)^3)$ flops overall.

In summary, the total complexity of the algorithm is

$$\mathcal{O}(dnr^3(r^* + p) + dnr^2(r^* + p)^2 + dn(r^* + p)^3)$$

Since $r^* + p \leq r^2$, we can simplify this expression to

$$\mathcal{O}(dnr^3(r^* + p) + dnr^2(r^* + p)^2)$$

Moreover, we may introduce the additional assumption that $r \leq r^* + p$ (i.e. that we do not wish to compress \mathcal{C} to TT-ranks which are lower than the ones of the two factors \mathcal{A} and \mathcal{B}) to conclude that the final cost is $\mathcal{O}(dnr^2(r^* + p)^2)$ flops.

Besides being cheaper than the deterministic approach (which has been analyzed in Section 5.2), Algorithm 14 also has lower memory requirements: indeed, just

$$\mathcal{O}(dr^2(r^* + p) + dn(r^* + p)^2)$$

memory is needed.

Given this procedure, we can easily tackle the rounding problem by modifying Algorithm 12: it suffices to replace the call to Algorithm 11 with one to Algorithm 14. This leads to Algorithm 15, where the integers m and q represent respectively the sampling step and the recompression gap (see Section 4.3.3).

A complexity analysis can be carried out similarly to the one of the original algorithm. For ease of notation we assume that both the initial and the final TT-ranks

Algorithm 15: Randomized rounding of the Hadamard product of two tensors in the TT format

Data: tensors $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{A}_1, \dots, \mathcal{A}_d$ and $\mathcal{B}_1, \dots, \mathcal{B}_d$ respectively, $\varepsilon > 0$, guess for the TT-ranks $(1, r_1^{(0)}, \dots, r_{d-1}^{(0)}, 1)$, positive integers q and m

Result: $\mathcal{C}^* \in \mathbb{R}^{n_1 \times \dots \times n_d}$ in the TT format with cores $\mathcal{C}_1^*, \dots, \mathcal{C}_d^*$, such that $\|\mathcal{C}^* - \mathcal{A} * \mathcal{B}\|_F \leq \varepsilon \|\mathcal{A} * \mathcal{B}\|_F$ holds with high probability

```

1  l := 0
2  do
3      C' := Algorithm 14(A, B, (1, r_1^{(l)}, ..., r_{d-1}^{(l)}, 1), 0)
4      C* := Algorithm 8(C', ε)                                     % from right to left
5      (1, r_1^*, ..., r_{d-1}^*, 1) = TT-ranks(C*)
6      changed := false
7      for j = 1, ..., d - 1 do
8          if r_j^{(l)} ≥ r_j^* + q then
9              r_j^{(l+1)} := r_j^{(l)}
10             else
11                 r_j^{(l+1)} := r_j^{(l)} + m
12                 changed := true
13             end
14         end
15         l := l + 1
16 while changed

```

are uniform:

$$r_1^{(0)} = \dots = r_{d-1}^{(0)} = r^{(0)} \quad \text{and} \quad \hat{r}_1 = \dots = \hat{r}_{d-1} = \hat{r}$$

Also, we assume that $\hat{r} = r^{(0)} + km - q$ for some integer k .

Thus, the l^{th} iteration (out of k) of Algorithm 15 takes $\mathcal{O}(dnr^2(r^{(0)} + lm)^2)$ flops because of Algorithm 14. This leads to an overall complexity of

$$\sum_{l=1}^k \mathcal{O}(dnr^2(r^{(0)} + lm)^2) = \mathcal{O}(dnr^2(r^{(0)} + km)^2 k) = \mathcal{O}(dnr^2(\hat{r} + q)^2 k)$$

5.3 Numerical examples

5.3.1 Square of a 2-dimensional rational function through a Hadamard power

We consider the matrix \mathbf{A} defined in Section 3.3.2, with $n = 10^5$ and $\gamma = 1$, and we compute its low-rank approximation with precision $\varepsilon = 10^{-8}$ in the Frobenius norm:

$$\mathbf{A} \approx \tilde{\mathbf{A}} = \mathbf{W}_A \mathbf{Z}_A^\top$$

In particular the rank of the decomposition is $r_A = 30$.

Now we wish to compute a low-rank approximation of $\tilde{\mathbf{A}} * \tilde{\mathbf{A}}$ with fixed precision $\varepsilon \in \{10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. To this aim, we apply the two approaches (deterministic and randomized) described in Section 5.1. In particular, we choose $q \in \{1, \dots, 6\}$ for the randomized algorithm. Due to the smoothness of $f_{2,2}(x, y) = (x + y)^{-2}$, the singular value decay of $\mathbf{A} * \mathbf{A}$ is still (at least) exponential. Hence, we expect small values of q to be sufficient to achieve the required precision.

We run the algorithms in the environment described in Section 3.3.1, and we plot the results in Figures 5.1 and 5.2.

In Figure 5.1 we show the relative errors of the final low-rank approximation. Similarly to the results of Section 3.3.3, we can observe that q plays a role analogous to the oversampling parameter in the fixed-rank case. In particular, we can verify that small values of q (e.g. $q = 3$) are enough to achieve the prescribed precision. Also, we report that the final rank of the dyadic approximation (after the recompression step in the randomized case) is always the same in both approaches, reaching the minimal size for which the low-rank approximation problem is feasible.

In Figure 5.2 we compare the computing time of the deterministic and randomized approaches. We can see that Algorithm 13 is always much faster. In particular, since the deterministic algorithm depends only weakly on ε , the relative efficiency of the randomized approach is more pronounced for larger values of ε .

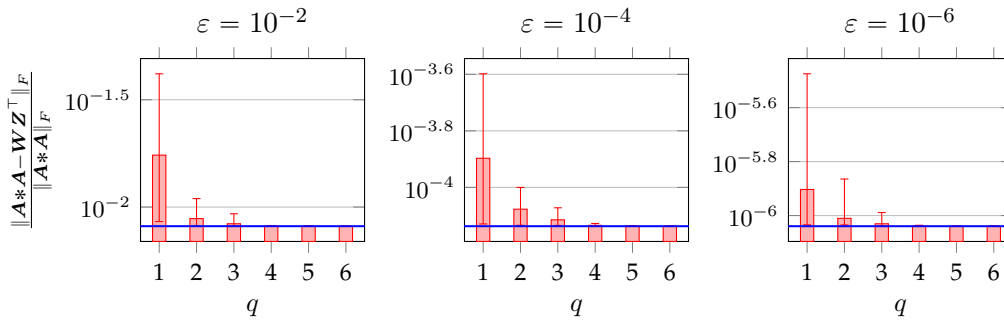


FIGURE 5.1: Relative error versus sampling parameter for Algorithm 13. The error bars indicate the range of values obtained in $N = 20$ simulations. In blue the result of the deterministic approach.

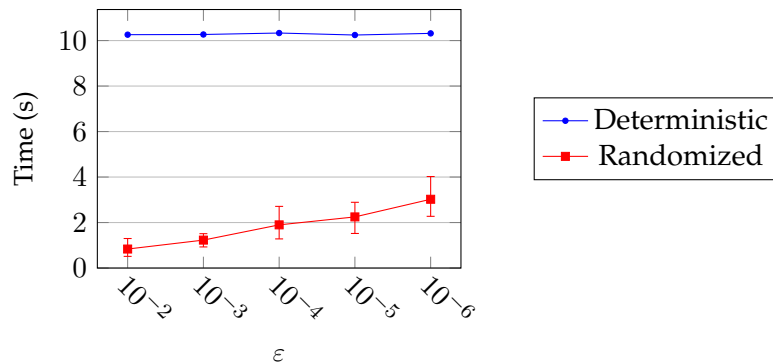


FIGURE 5.2: Computing time for Algorithm 13 and the deterministic counterpart. For the randomized algorithm, we choose a sampling parameter $q = 3$. The error bars indicate the range of values obtained in $N = 20$ simulations.

5.3.2 Hadamard product of high-dimensional rational functions

Let \mathcal{A} be the tensor defined in Section 4.4.2, with $d = 15$, $n = 30$, and $\gamma = 1$. Moreover, we consider \mathcal{B} , the tensor obtained in the same way but with $\gamma = 1/2$. Such tensors are generated by applying an exponential sum approximation with 20 terms. With this choice, we obtain tensors with uniform TT-ranks $(1, 20, \dots, 20, 1)$, which provide approximate discretizations of

$$f_{d,1}(x_1, \dots, x_d) = (x_1 + \dots + x_d)^{-1} \quad \text{and} \quad f_{d,1/2}(i_1, \dots, i_d) = (x_1 + \dots + x_d)^{-1/2}$$

We wish to recompress the Hadamard product $\mathcal{C} = \mathcal{A} * \mathcal{B}$ down to precision $\varepsilon \in \{10^{-2}, 10^{-4}, 10^{-6}, 10^{-8}, 10^{-10}\}$, using two approaches: the first is the explicit computation of \mathcal{C} , followed by the execution of Algorithm 8; the second is Algorithm 15 with $m = 2$ and $q = 1$, without any initial guess for the compressed TT-ranks (i.e. we set $r_1^{(0)} = \dots = r_{d-1}^{(0)} = 0$). Such small choices for m and q are motivated by the fact that the spectra of all unfoldings of \mathcal{C} are (at least) exponentially decaying, due to the regularity properties of $f_{d,3/2}$ in $(0, \frac{n+1}{10})^d$.

Figures 5.3 and 5.4 show the results of the simulations, obtained in the environment described in Section 3.3.1. From Figure 5.3 (left) we can observe that the two approaches yield similar errors, achieving the desired accuracy in almost all simulations. However, even in case of failure, the value of the relative error is within a small factor of ε .

The computing time is shown in Figure 5.3 (right), where we can see that the randomized approach is much faster than the deterministic one. Also, it is interesting to note that, for $\varepsilon \geq 10^{-6}$, it is faster to run Algorithm 15 than to carry out just the preliminary computation of \mathcal{C} , which is required for the deterministic approach.

In Figure 5.4 we show the maximum TT-ranks at the end of the two algorithms. In particular, for the randomized approach we also show the TT-ranks before the last recompression step, i.e. the final TT-ranks of \mathcal{C}' in Algorithm 15. We can observe that a small amount of oversampling (measured as the variation between the TT-ranks before and after recompression) is enough to guarantee the correctness of the randomized procedure, which (except for a single sample with $\varepsilon = 10^{-4}$) achieves the same TT-ranks as the deterministic approach.

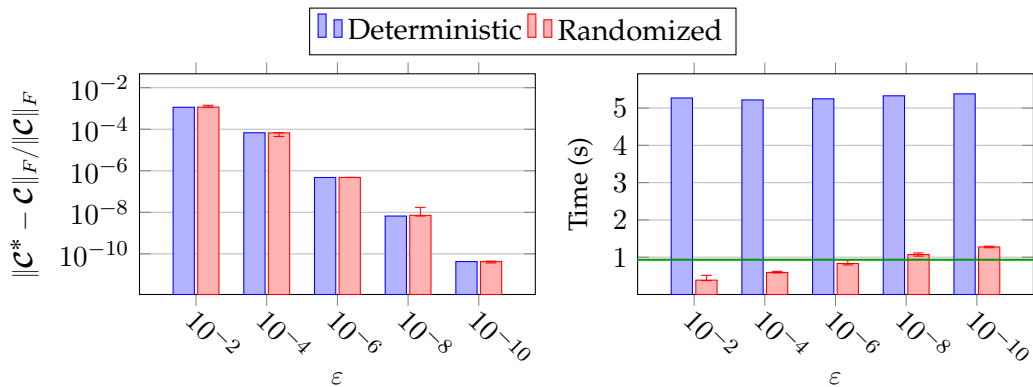


FIGURE 5.3: Relative error (left) and computing time (right) for the deterministic approach and for Algorithm 15. The bars connect the minimum and the maximum values obtained over $N = 20$ simulations. The green line in the right plot represents the time needed for the computation of \mathcal{C} in the deterministic approach.

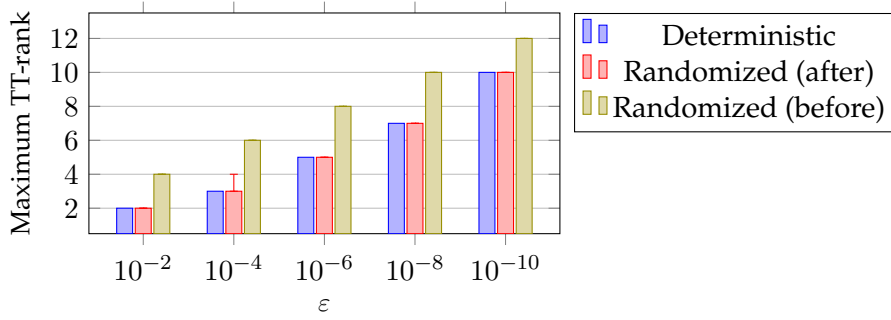


FIGURE 5.4: Maximum TT-rank of the compressed tensor for the deterministic approach and for Algorithm 15. For the randomized algorithm the plot shows the modal values of the maximal TT-ranks over $N = 20$ simulations, both before and after recompression. Also, the bars show the range of values obtained.

5.3.3 Hadamard product of random tensors

We build the random tensors \mathcal{A} and \mathcal{B} by following the procedure described in Section 4.4.3. In particular, we choose $n = 50$, $d \in \{10, 20, 40\}$, and $r \in \{4, 6, \dots, 40\}$. For each value of r , we choose an exponentially decaying spectrum $\sigma_j = a^{1-j}$ (for $j \in \{1, \dots, r\}$), with $a > 1$ such that σ_r is approximately equal to the machine precision, i.e. we set $a \simeq \varepsilon_m^{1/(1-r)}$.

We want to recompress $\mathcal{C} = \mathcal{A} * \mathcal{B}$ down to uniform TT-ranks $(1, r, \dots, r, 1)$. To this aim, we follow several strategies:

- we compute \mathcal{C} and apply Algorithm 7,
- we compute \mathcal{C} and apply Algorithm 11 (with $p = 5$),
- we apply Algorithm 14 (with $p = 5$).

In particular, we want to compare the computing times of the three approaches. The results are shown in Figure 5.5 for the computational environment described in Section 3.3.1.

We can observe that, for all choices of d , the three approaches have similar speeds for small values of r : actually, Algorithm 14 is outperformed by the other two procedures for $r \leq 8$. The fact that the deterministic approach is more efficient for small tensors appears reasonable, since it relies more heavily on highly optimized built-in MATLAB[®] functions.

In contrast, for larger values of r , both approaches which are based on the explicit computation of \mathcal{C} perform quite poorly when compared to Algorithm 14. For instance, if we consider $d = 40$ and $r = 30$, Algorithm 14 has an average running time of 6.28 seconds, versus 35.99 seconds for Algorithm 11 (speed-up ratio: 5.73) and 205.3 seconds for Algorithm 7 (speed-up ratio: 32.7).

In particular, as in the previous example, we can identify several configurations ($r \geq 22$) where Algorithm 14 is faster than the computation of \mathcal{C} , which represents a preliminary step in the other approaches.

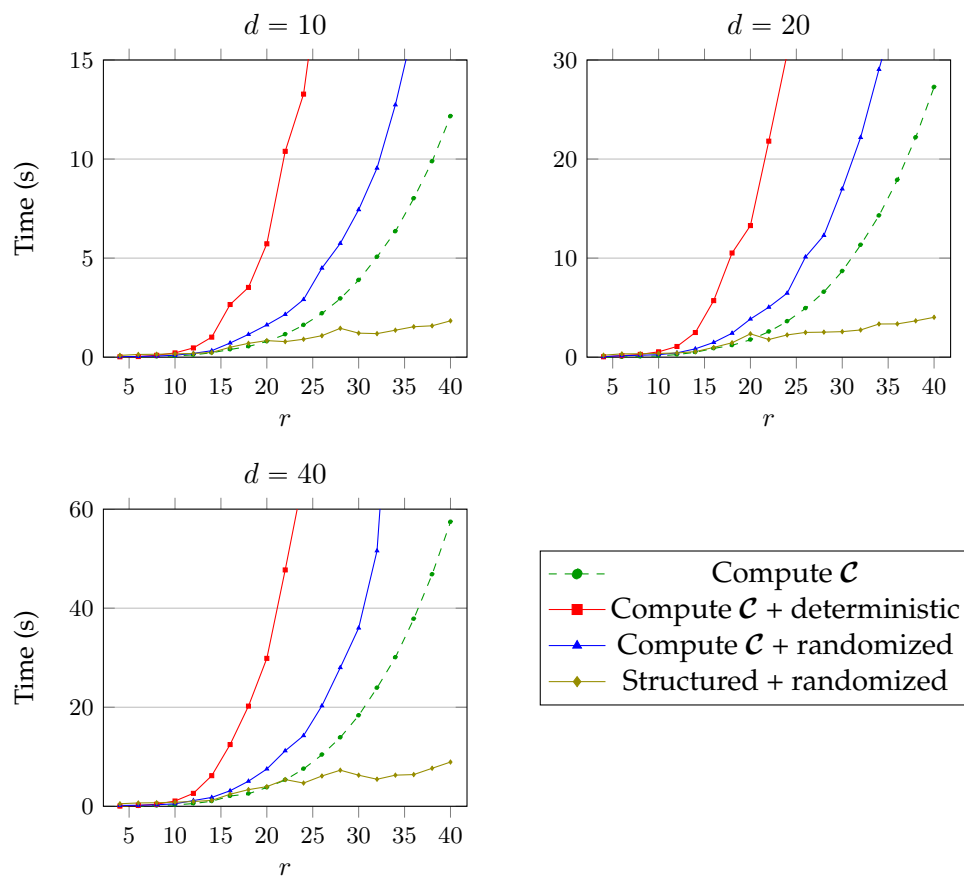


FIGURE 5.5: Computing time versus initial TT-ranks for Algorithms 7, 11 and 14. For the two randomized approaches we show the average value over $N = 20$ simulations. In green the time needed just for the computation of \mathcal{C} .

Chapter 6

Conclusion

We have shown that randomized approaches to low-rank matrix and tensor approximation are quite satisfactory. Indeed, at the cost of a small oversampling, they can achieve a similar accuracy as the corresponding deterministic techniques.

Moreover, they are often more efficient than any deterministic algorithm. This is mainly due to the high degree of flexibility allowed by randomness, which leads to a better exploitation of the structure of the data which needs to be compressed.

For instance, in the compression of tensors in the TT format, we have described how randomness allows to skip the orthogonalization step, which is required in any (efficient) deterministic approach. A similar effect can be observed in the case of Hadamard products, where we have derived randomized approaches which avoid the computation of the full products.

Among the several topics which have been considered in this thesis, the most striking results are the ones obtained on tensors in the TT format. Indeed, Algorithms 11 and 12 (in the standard case) and Algorithms 14 and 15 (in the Hadamard product case) represent very interesting alternatives to their deterministic counterpart, both in terms of efficiency and memory requirements.

Moreover, it is particularly impressive to observe the accuracy achieved by all algorithms for the compression of tensors in the TT format, considering that they rely on “rank-1” random vectors. For such vectors, a thorough theoretical analysis, generalizing the available results from the Gaussian case, is still missing, even though a relevant effort is being made in current research. Still, we have shown with several numerical examples that “rank-1” random vectors provide an efficient way to obtain samples from the ranges of large matrices, for which other techniques would be too expensive.

Bibliography

- [Aff+87] I. Affleck et al. “Rigorous results on valence-bond ground states in antiferromagnets”. In: *Phys. Rev. Lett.* 59 (7 1987), pp. 799–802. DOI: 10.1103/PhysRevLett.59.799.
- [BH05] D. Braess and W. Hackbusch. “Approximation of $1/x$ by exponential sums in $[1, \infty)$ ”. In: *IMA Journal of Numerical Analysis* 25.4 (2005). DOI: 10.1093/imanum/dri015.
- [Bis09] C.M. Bishop. *Pattern recognition and machine learning*. Information Science and Statistics. Springer, 2009. ISBN: 978-0387-31073-2.
- [BNT07] I. Babuška, F. Nobile, and R. Tempone. “A Stochastic Collocation Method for Elliptic Partial Differential Equations with Random Input Data”. In: *SIAM Journal on Numerical Analysis* 45.3 (2007), pp. 1005–1034. DOI: 10.1137/050645142.
- [CD05] Z. Chen and J.J. Dongarra. “Condition Numbers of Gaussian Random Matrices”. In: *SIAM Journal on Matrix Analysis and Applications* 27.3 (2005), pp. 603–620. DOI: 10.1137/040616413.
- [Cha87] T.F. Chan. “Rank revealing QR factorizations”. In: *Linear Algebra and its Applications* 88 (1987), pp. 67–82. DOI: 10.1016/0024-3795(87)90103-0.
- [Dix83] J.D. Dixon. “Estimating Extremal Eigenvalues and Condition Numbers of Matrices”. In: *SIAM Journal on Numerical Analysis* 20.4 (1983), pp. 812–814. DOI: 10.1137/0720053.
- [EY36] C. Eckart and G. Young. “The approximation of one matrix by another of lower rank”. In: *Psychometrika* 1.3 (1936), pp. 211–218. ISSN: 1860-0980. DOI: 10.1007/BF02288367.
- [FNW92] M. Fannes, B. Nachtergaele, and R.F. Werner. “Finitely correlated states on quantum spin chains”. In: *Communications in Mathematical Physics* 144.3 (1992), pp. 443–490. DOI: 10.1007/BF02099178.
- [GE96] M. Gu and S.C. Eisenstat. “Efficient Algorithms for Computing a Strong Rank-Revealing QR Factorization”. In: *SIAM Journal on Scientific Computing* 17.4 (1996), pp. 848–869. DOI: 10.1137/0917055.
- [GL17] M. Griebel and G. Li. “On the decay rate of the singular values of bivariate functions”. In: 2017. URL: <http://wissrech.ins.uni-bonn.de/research/pub/li/INSPreprint1702.pdf>.
- [GL89] G.H. Golub and C.F. van Loan. *Matrix Computations*. Johns Hopkins Series in the Mathematical Sciences 3. The Johns Hopkins University Press, 1989.

- [HMT11] N. Halko, P.G. Martinsson, and A.J. Tropp. "Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions". In: *SIAM Review* 53.2 (2011), pp. 217–288. DOI: 10.1137/090771806.
- [HP92] Y.P. Hong and C.-T. Pan. "Rank-Revealing QR Factorizations and the Singular Value Decomposition". In: *Mathematics of Computation* 58.197 (1992), pp. 213–232. ISSN: 00255718, 10886842.
- [KB09] T.G. Kolda and B.W. Bader. "Tensor Decompositions and Applications". In: *SIAM Review* 51.3 (2009), pp. 455–500. DOI: 10.1137/07070111X.
- [KP16] D. Kressner and L. Periša. *Recompression of Hadamard Products of Tensors in Tucker Format*. Tech. rep. 2016. URL: http://sma.epfl.ch/~anchpcommon/publications/ttensors_pp.pdf.
- [Lib+07] E. Liberty et al. "Randomized algorithms for the low-rank approximation of matrices". In: *Proceedings of the National Academy of Sciences* 104.51 (2007), pp. 20167–20172. DOI: 10.1073/pnas.0709640104.
- [Lub08] C. Lubich. *From Quantum to Classical Molecular Dynamics: Reduced Models and Numerical Analysis*. Zurich Lectures in Advanced Mathematics. European Mathematical Society Publishing House, 2008. DOI: 10.4171/067.
- [McL16] W. McLean. "Exponential sum approximations for $t^{-\beta}$ ". In: *ArXiv e-prints* (2016). arXiv: 1606.00123.
- [Mir60] L. Mirsky. "Symmetric gauge functions and unitarily invariant norms". In: *The Quarterly Journal of Mathematics* 11.1 (1960), pp. 50–59. DOI: 10.1093/qmath/11.1.50.
- [Oru14] R. Orus. "A Practical Introduction to Tensor Networks: Matrix Product States and Projected Entangled Pair States". In: *Annals Phys.* 349 (2014), pp. 117–158. DOI: 10.1016/j.aop.2014.06.013.
- [Ose11] I.V. Oseledets. "Tensor-Train Decomposition". In: *SIAM Journal on Scientific Computing* 33.5 (2011), pp. 2295–2317. DOI: 10.1137/090752286.
- [OT09] I.V. Oseledets and E.E. Tyrtyshnikov. "Breaking the Curse of Dimensionality, Or How to Use SVD in Many Dimensions". In: *SIAM Journal on Scientific Computing* 31.5 (2009), pp. 3744–3759. DOI: 10.1137/090748330.
- [Par98] B. Parlett. *The Symmetric Eigenvalue Problem*. Society for Industrial and Applied Mathematics, 1998. DOI: 10.1137/1.9781611971163.
- [PP07] J. Persson and L. von Persson. "Pricing European multi-asset options using a space-time adaptive FD-method". In: *Computing and Visualization in Science* 10.4 (2007), pp. 173–183. DOI: 10.1007/s00791-007-0072-y.
- [RW05] C.E. Rasmussen and C.K.I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005. ISBN: 026218253X.
- [Sch04] K. Schäcke. "On the Kronecker Product". MA thesis. University of Waterloo, 2004.

- [Ste16] M. Steinlechner. “Riemannian Optimization for Solving High-Dimensional Problems with Low-Rank Tensor Structure”. PhD thesis. EPFL, 2016. DOI: doi:10.5075/epfl-thesis-6958.
- [Ven+07] O. Vendrell et al. “Full-dimensional (15-dimensional) quantum-dynamical simulation of the protonated water dimer. I. Hamiltonian setup and analysis of the ground vibrational state”. In: *The Journal of Chemical Physics* 127.18 (2007), p. 184302. DOI: 10.1063/1.2787588.
- [ZJD15] K. Zhong, P. Jain, and I.S. Dhillon. “Efficient Matrix Sensing Using Rank-1 Gaussian Measurements”. In: *Algorithmic Learning Theory: 26th International Conference, ALT 2015, Banff, AB, Canada, October 4-6, 2015, Proceedings*. Springer International Publishing, 2015, pp. 3–18. DOI: 10.1007/978-3-319-24486-0_1.