

# Error Analysis of TT-format Tensor Algorithms

Dario Fasino and Eugene E. Tyrtyshnikov

**Abstract** The tensor train (TT) decomposition is a representation technique for arbitrary tensors, which allows efficient storage and computations. For a  $d$ -dimensional tensor with  $d \geq 2$ , that decomposition consists of two ordinary matrices and  $d - 2$  third-order tensors. In this paper we prove that the TT decomposition of an arbitrary tensor can be computed (or approximated, for data compression purposes) by means of a backward stable algorithm based on computations with Householder matrices. Moreover, multilinear forms with tensors represented in TT format can be computed efficiently with a small backward error.

## 1 Introduction

The *tensor train decomposition* is a representation technique which allows compact storage and efficient computations with arbitrary tensors. The origins of this representation can be traced back to a brief paper by I. Oseledets and E. Tyrtyshnikov dating 2009 [8], while its popularization is mainly due to the subsequent papers [6, 9]. Nowadays, the tensor train decomposition is a computationally powerful tool that offers viable and convenient alternatives to classical (e.g., Tucker, CP) tensor representations [1, 4], in particular for the approximation of solutions of high dimensional problems. As shown in, e.g., [5, 7], certain computations with large-scale structured matrices and vectors can be conveniently recast in terms of tensor train representations.

---

Dario Fasino  
Department of Mathematics, Computer Science and Physics, University of Udine, Udine, Italy,  
e-mail: dario.fasino@uniud.it

Eugene E. Tyrtyshnikov  
Institute of Numerical Mathematics of Russian Academy of Sciences, Russia, e-mail: eugene.tyrtyshnikov@gmail.com

Basically, a tensor train (TT) decomposition of a  $d$ -dimensional tensor  $A$  with size  $n_1 \times n_2 \times \dots \times n_d$  is a sequence  $G_1, \dots, G_d$  of tensors of order 2 or 3; the size of  $G_i$  is  $r_{i-1} \times n_i \times r_i$  with  $r_0 = r_d = 1$  (that is,  $G_1$  and  $G_d$  are ordinary matrices) and

$$A(i_1, i_2, \dots, i_d) = \sum_{j_1=1}^{r_1} \dots \sum_{j_{d-1}=1}^{r_{d-1}} G_1(i_1, j_1) G_2(j_1, i_2, j_2) \dots G_d(j_{d-1}, i_d). \quad (1)$$

We will denote by  $\text{TT}(G_1, \dots, G_d)$  the tensor identified by the right hand side of (1).

In this paper, we present a backward error analysis of two algorithms, originally devised in [9], which perform computations with tensors in TT-format. The first algorithm produces an exact or approximate TT decomposition  $G_1, \dots, G_d$  of an arbitrary  $d$ -dimensional tensor  $A$  given in functional form, depending on a tolerance parameter  $\varepsilon$ . If  $\varepsilon = 0$  then the output of the algorithm is an exact TT decomposition, that is,  $A = \text{TT}(G_1, \dots, G_d)$ . If  $\varepsilon > 0$  then  $\text{TT}(G_1, \dots, G_d)$  is an  $\mathcal{O}(\varepsilon)$ -approximation of  $A$  which can realize significant savings in memory space. The computational core of the algorithm is a suitable (approximate) matrix factorization that, in the original paper, relies on SVD computations. We prove that analogous performances and backward stability can be obtained by means of QR factorizations based on Householder transformations.

The second algorithm computes the *contraction* (i.e., multilinear form) of a given  $d$ -dimensional tensor  $A$  and vectors  $v^{(1)}, \dots, v^{(d)}$ ,

$$\alpha = \sum_{i_1=1}^{n_1} \dots \sum_{i_d=1}^{n_d} A(i_1, i_2, \dots, i_d) v_{i_1}^{(1)} \dots v_{i_d}^{(d)}, \quad (2)$$

where  $A$  is known in TT format. By means of known error bounds for inner products in floating point arithmetic [3], we prove backward stability of the proposed algorithm under very general hypotheses on the evaluation order of the summations. More precisely, if  $A = \text{TT}(G_1, \dots, G_d)$  and no underflows or overflows are encountered then the output computed by the algorithm in floating point arithmetic is the exact contraction of  $\hat{A} = \text{TT}(G_1 + \Delta G_1, \dots, G_d + \Delta G_d)$  and  $v^{(1)}, \dots, v^{(d)}$  where  $|\Delta G_i| \leq (n_i + r_{i-1})u|G_i| + \mathcal{O}(u^2)$  and  $u$  is the machine precision.

After setting up some basic notations and concepts, in Section 3 we present the algorithm for computing the TT-representation of a tensor and analyze its numerical stability. Next, we discuss in Section 4 the computation in computer arithmetic of the multilinear form (2) with a tensor in TT-format. A final appendix contains a complementary result.

## 2 Notations and Preliminaries

We refer to [1, Ch. 12] and [4] for notations and fundamental concepts on tensors and basic multilinear operations. Vectors and matrices are denoted by lower-case and upper-case italic letters, respectively, and higher order tensors by upper-case

sans-serif letters:  $x$ ,  $X$ , and  $\mathbb{X}$ . A tensor  $\mathbb{X}$  of order  $d$  is a multiway array of size  $n_1 \times n_2 \times \cdots \times n_d$ , where  $n_k$  is the size of the  $k$ th dimension or mode. A vector is a first-order tensor, and a matrix is a second-order tensor. The  $(i_1, i_2, \dots, i_d)$ th entry of  $\mathbb{X} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$  is denoted by  $X_{i_1, i_2, \dots, i_d}$  or, alternatively,  $\mathbb{X}(i_1, i_2, \dots, i_d)$  to avoid long, multiple subscripts. Lower case greek letters denote real numbers.

Throughout this paper, the prototypical tensor is a  $d$ -dimensional array  $A$  with dimensions  $n_1 \times n_2 \times \cdots \times n_d$ . In particular, the scalar  $d$  is used exclusively for the order of  $A$ , and the scalars  $n_1, n_2, \dots, n_d$  are reserved for  $A$ 's dimensions. The *size* of  $A$  is the number of entries,  $N(A) = n_1 n_2 \cdots n_d$ . The *vectorization* of  $A$  is the vector  $\text{vec}(A) \in \mathbb{R}^{N(A)}$  whose  $i$ th entry is the  $i$ th entry of  $A$  according to the lexicographic ordering of the indices. The Matlab-style function `reshape` is defined in terms of the vectorization operator as follows. Let  $m_1, m_2, \dots, m_k$  be integers such that  $N(A) = m_1 m_2 \cdots m_k$ . Then,

$$B = \text{reshape}(A, [m_1, m_2, \dots, m_k])$$

is the tensor  $B \in \mathbb{R}^{m_1 \times m_2 \times \cdots \times m_k}$  such that  $\text{vec}(A) = \text{vec}(B)$ . In particular, for  $k = 1, \dots, d-1$ ,

$$A_k = \text{reshape}(A, [\prod_{i=1}^k n_i, \prod_{i=k+1}^d n_i])$$

is the  $k$ -th *unfolding matrix* of  $A$ .

The  $k$ -mode *product* of a tensor  $A \in \mathbb{R}^{n_1 \times \cdots \times n_d}$  by a matrix  $M \in \mathbb{R}^{n_k \times m}$ , denoted by  $A \times_k M$ , is an  $(n_1 \times \cdots \times n_{k-1} \times m \times n_{k+1} \times \cdots \times n_d)$ -tensor of which the entries are given by

$$(A \times_k M)(i_1, \dots, i_d) = \sum_{j=1}^{n_k} A(i_1, \dots, i_{k-1}, j, i_{k+1}, \dots, i_d) M_{j, i_k}.$$

The latter definition extends trivially to the case where  $M$  is a vector, by treating it as a  $n_k \times 1$  matrix. The  $k$ -mode product satisfies the property

$$(A \times_i B) \times_j C = (A \times_j C) \times_i B, \quad (3)$$

so that notations like  $A \times_i B \times_j C$  can be used without ambiguity.

If  $A$  is a vector, a matrix or a tensor, we denote by  $|A|$  the componentwise absolute value of  $A$ . Inequalities between tensors hold componentwise. Finally, the Frobenius inner product of two matrices  $A, B \in \mathbb{R}^{m \times n}$  is  $\langle A, B \rangle = \text{trace}(A^T B) = \text{vec}(A)^T \text{vec}(B)$ , and the associated matrix norm is  $\|A\|_F = \|\text{vec}(A)\|_2$ . The latter is extended in an obvious way to arbitrary tensors.

## 2.1 The Tensor Train Format for Multidimensional Arrays

A tensor train decomposition of  $A \in \mathbb{R}^{n_1 \times \cdots \times n_d}$  is a sequence  $G_1, \dots, G_d$  of tensors, called *carriages* (or *cores* [6, 7]), such that the size of  $G_i$  is  $r_{i-1} \times n_i \times r_i$  with  $r_0 =$

$r_d = 1$  (that is,  $G_1$  and  $G_d$  can be understood as ordinary matrices) and fulfilling the identity (1), that is,  $A = \text{TT}(G_1, \dots, G_d)$ .

*Example 1.* Let  $A$  be a  $10 \times 20 \times 30 \times 40$  tensor with TT-ranks 5, 6, 7. The TT-format of  $A$  is made of 4 carriages, whose dimensions are as follows:

$$\begin{aligned} G_1 &: 1 \times 10 \times 5 \\ G_2 &: 5 \times 20 \times 6 \\ G_3 &: 6 \times 30 \times 7 \\ G_4 &: 7 \times 40 \times 1. \end{aligned}$$

An alternative viewpoint on the decomposition (1) is

$$A(i_1, i_2, \dots, i_d) = G'_1(i_1)G'_2(i_2) \cdots G'_d(i_d),$$

where now  $G'_k(i_k)$  is an  $r_{k-1} \times r_k$  matrix depending on the integer parameter  $i_k$ . The numbers  $r_1, \dots, r_{d-1}$  are called *TT-ranks* of  $A$ . As shown in, e.g., [6, 9],  $r_k$  is bounded from below by  $\text{rank}(A_k)$ . Moreover, a TT decomposition with  $r_k = \text{rank}(A_k)$  always exists and can be computed by the algorithm shown in [9], which is recalled in the next section. It is often the case that an exact or approximate TT-format of a given tensor yields considerable savings in terms of memory space with respect to other representation techniques.

*Remark 1.* It is sometimes convenient to assume that  $G_1$  and  $G_d$  are not three-dimensional but two-dimensional with sizes  $n_1 \times r_1$  and  $r_{d-1} \times n_d$ , respectively. For that reason, in what follows we will use indifferently the notations  $G_1$  and  $G_1$  to denote the first carriage.

### 3 Full-to-TT Compression

In this section we address backward stability properties of the compression algorithm from [9], which is recalled hereafter as Algorithm 3.1. This algorithm produces an exact or approximate TT decomposition  $G_1, \dots, G_d$  of a given tensor  $A$  with  $d \geq 2$ , depending on a tolerance  $\varepsilon$ . If  $\varepsilon = 0$  then the output of the algorithm is an exact TT decomposition, that is  $A = \text{TT}(G_1, \dots, G_d)$  with  $r_k = \text{rank}(A_k)$  for  $k = 1, \dots, d-1$ , see [9, Thm. 2.1]. If  $\varepsilon > 0$  then  $r_k \leq \text{rank}(A_k)$  and  $\text{TT}(G_1, \dots, G_d)$  is an  $\mathcal{O}(\varepsilon)$ -approximation of  $A$  which can realize significant savings in memory space with respect to the exact TT-decomposition.

As a basic compression device we suppose to have at our disposal a generic, black-box algorithm having the interface

$$[M, N, r] = \text{compress}(A, \varepsilon) \tag{4}$$

which, for any matrix  $A \in \mathbb{R}^{m \times n}$  and tolerance  $\varepsilon \geq 0$  returns an integer  $r$  and two matrices  $M \in \mathbb{R}^{m \times r}$  and  $N \in \mathbb{R}^{r \times n}$  such that

$$\text{rank}(M) = \text{rank}(N) = r, \quad \|MN - A\|_F \leq \varepsilon \|A\|_F. \quad (5)$$

In particular, if  $\varepsilon = 0$  then  $A = MN$  is a rank decomposition of  $A$ . In [9] this algorithm is realized by means of a truncated SVD of  $A$ . In that case, one has

$$\|A - MN\|_F = \min_{\text{rank}(X) \leq r} \|A - X\|_F, \quad r = \min_{\|A - X\|_F \leq \varepsilon \|A\|_F} \text{rank}(X).$$

Consequently, if  $\varepsilon > 0$  then TT-ranks  $r_1, \dots, r_{d-1}$  computed by the resulting procedure are in some sense optimal. However, other rank-revealing factorizations can be usefully adopted for the purpose of computing the TT-format of (an approximation of) the given tensor.

Observe that Algorithm 3.1 is based on a finite iteration. Each iteration is entirely based on matrix computations. In particular, if the argument  $A$  is a matrix then the output consists of the two matrices computed by `compress`. In the pseudocode here below, the intermediate matrices  $B_1, \dots, B_{d-1}$  are introduced to improve readability; in a more implementation-oriented description, these variables can share the same name and memory space.

---

**Algorithm 3.1** Full-to-TT compression, iterative version

---

**Input:** tensor  $A$  of size  $n_1 \times n_2 \times \dots \times n_d$  and local accuracy bound  $\varepsilon$

**Output:** tensor carriages  $G_1, \dots, G_d$

```

1: function [ $G_1, \dots, G_d$ ] = FULL_TO_TT( $A, \varepsilon$ )
2:    $n := N(A)$ 
3:    $A_1 := \text{reshape}(A, [n_1, n/n_1])$ 
4:   [ $G_1, B_1, r_1$ ] := compress( $A_1, \varepsilon$ )
5:   for  $k = 2 \dots d - 1$  do
6:      $n_R := n_R / n_k$ 
7:      $B_{k-1} := \text{reshape}(B_{k-1}, [r_{k-1} n_k, n_R])$ 
8:     [ $G_k, B_k, r_k$ ] := compress( $B_{k-1}, \varepsilon$ )
9:      $G_k := \text{reshape}(G_k, [r_{k-1}, n_k, r_k])$ 
10:  end for
11:   $G_d := B_{d-1}$ 
12: end function

```

---

*Example 2.* Suppose that Algorithm 3.1 is applied to the fourth order tensor  $A$  in Example 1. Then Algorithm 3.1 performs three for-loops; the dimensions of the matrices  $B_1$  and  $B_2$  before and after being reshaped in step 7 are as follows:

$B_1$  :  $5 \times 24000$  and  $100 \times 1200$

$B_2$  :  $6 \times 1200$  and  $180 \times 40$

At the third iteration the matrix  $B_3$  is  $7 \times 40$  and is renamed  $G_4$  without being reshaped.

Algorithm 3.2 is a recursive version of Algorithm 3.1, namely, the two algorithms compute the same function, but Algorithm 3.2 calls itself internally. Termination is ensured by the fact that the order of the tensor argument diminishes by one at each

recursive call. In fact, if  $d = 2$  then the sought decomposition is obtained immediately by one call to `compress`, as in the iterative version. When  $d > 2$ , the matrix  $B_1$  computed in step 4 together with the carriage  $G_1$  is the first unfolding matrix of the  $(d - 1)$ -dimensional tensor  $B$ , whose TT decomposition is computed by the recursive call in step 9. In that step the carriage  $G_2$  is obtained as a matrix of dimension  $(r_1 n_2) \times r_2$ . In the subsequent steps that matrix is reshaped into a tensor of dimension  $r_1 \times n_2 \times r_2$  and the computation is complete.

---

**Algorithm 3.2** Full-to-TT compression, recursive version

---

**Input:** tensor  $A$  of size  $n_1 \times n_2 \times \dots \times n_d$  and local accuracy bound  $\varepsilon$

**Output:** tensor carriages  $G_1, \dots, G_d$

```

1: function  $[G_1, \dots, G_d] = \text{Full\_to\_TT}(A, \varepsilon)$ 
2:    $n := N(A)$ 
3:    $A_1 := \text{reshape}(A, [n_1, n/n_1])$ 
4:    $[G_1, B_1, r_1] := \text{compress}(A_1, \varepsilon)$ 
5:   if  $d = 2$  then
6:      $G_2 := B_1$ 
7:   else
8:      $B := \text{reshape}(B_1, [r_1 n_2, n_3, \dots, n_d])$  ▷ the order of  $B$  is  $d - 1$ 
9:      $[G_2, \dots, G_d] := \text{Full\_to\_TT}(B, \varepsilon)$  ▷ recursive call
10:     $r_2 := N(G_2)/(r_1 n_2)$ 
11:     $G_2 := \text{reshape}(G_2, [r_1, n_2, r_2])$  ▷  $G_2$  is tensorized
12:   end if
13: end function

```

---

*Remark 2.* Apart of reshaping, the matrix  $B_k$  computed in Algorithm 3.1 coincides with the matrix  $B_1$  computed in the  $(k - 1)$ -th recursive call of Algorithm 3.2.

In the subsequent analysis we assume that, in exact arithmetic, the function `compress` satisfies the following property: For any input matrix  $A$ , the output matrices  $M$  and  $N$  fulfil the identities

$$M^T M = I, \quad M^T (A - MN) = O. \quad (6)$$

Equations (5) and (6) can be met if `compress` is obtained from a truncated SVD as in [9] or from a (truncated) QR factorization. Indeed, the equations in (6) are fulfilled if and only if the matrix  $A$  admits a factorization

$$A = (M M') \begin{pmatrix} N \\ N' \end{pmatrix}$$

where the matrix  $Q = (M M')$  has orthonormal columns and  $N$  has full rank, in which case we have  $\|A - MN\| = \|M'N'\| = \|N'\|$  and  $\|N\| = \|M^T A\| \leq \|A\|$  in any unitarily invariant norm. The sufficiency of that condition is obvious. To prove necessity, consider the residual matrix  $R = A - MN$  and let  $R = M'N'$  be a factorization where the columns of  $M'$  are an orthonormal basis for the column space of  $R$  and  $N'$

has full rank. Since  $M^T M' N' = O$  the columns of  $M'$  must belong to the kernel of  $M^T$ , hence the matrix  $Q = (M M')$  has orthonormal columns.

The following theorem is a reworking of Theorem 2.2 from [9] to emphasize the role of the hypotheses placed on `compress`. A shorter proof is included for later reference.

**Theorem 1.** *Suppose that the function `compress` in (4) fulfills the conditions (5) and (6). Let  $\mathbb{T} = \text{TT}(G_1, \dots, G_d)$  where the tensor carriages  $G_1, \dots, G_d$  are computed from Algorithm 3.2 in exact arithmetic. Then  $\|A - \mathbb{T}\|_{\text{F}} \leq \varepsilon \sqrt{d-1} \|A\|_{\text{F}}$ .*

*Proof.* We proceed by induction on  $d$ . When  $d = 2$  the claim follows immediately from (5). For  $d > 2$ , consider the matrices  $G_1$  and  $B_1$  computed in step 4. By construction, the first unfolding matrix of  $\mathbb{T}$  admits the factorization  $T_1 = G_1 U_1$ , for some  $r_1 \times (n_2 \cdots n_d)$  matrix  $U_1$ . Since  $G_1^T (A_1 - G_1 B_1) = O$  and  $G_1^T G_1 = I$  by hypothesis, we have

$$\begin{aligned} \|A - \mathbb{T}\|_{\text{F}}^2 &= \|A_1 - G_1 B_1 + G_1 B_1 - G_1 U_1\|_{\text{F}}^2 \\ &= \|A_1 - G_1 B_1\|_{\text{F}}^2 + \|G_1 (B_1 - U_1)\|_{\text{F}}^2 + 2\langle G_1^T (A_1 - G_1 B_1), B_1 - U_1 \rangle \\ &= \|A_1 - G_1 B_1\|_{\text{F}}^2 + \|B_1 - U_1\|_{\text{F}}^2. \end{aligned}$$

Consider the  $(d-1)$ -dimensional tensor

$$\mathbb{U} = \text{reshape}(U_1, [r_1 n_2, n_3, \dots, n_d]),$$

whose first unfolding matrix is  $U_1$ . Then,

$$\|A - \mathbb{T}\|_{\text{F}}^2 \leq \varepsilon^2 \|A\|_{\text{F}}^2 + \|B - \mathbb{U}\|_{\text{F}}^2,$$

where  $B$  is obtained in step 8 of Algorithm 3.2. By construction  $\mathbb{U} = \text{TT}(G_2, \dots, G_d)$ , that is,  $\mathbb{U}$  is the tensor whose exact TT-decomposition is given by the tensor carriages obtained from Algorithm 3.2 with input  $B$ . By inductive hypothesis we have

$$\|B_1 - U_1\|_{\text{F}}^2 = \|B - \mathbb{U}\|_{\text{F}}^2 \leq (d-2)\varepsilon^2 \|B\|_{\text{F}}^2.$$

Moreover,

$$\|B\|_{\text{F}} = \|B_1\|_{\text{F}} = \|G_1^T A_1\|_{\text{F}} \leq \|A_1\|_{\text{F}} = \|A\|_{\text{F}},$$

and the claim follows.  $\square$

*Remark 3.* It is not difficult to prove that the hypothesis  $M^T M = I$  in (6) can be replaced by  $\|M\|_2 \|M^+\|_2 = 1$ , where  $M^+$  is the Moore–Penrose inverse of  $M$ , without affecting the claim of Theorem 1. In the proof, one has simply to replace  $\|G_1 (B_1 - U_1)\|_{\text{F}} = \|B_1 - U_1\|_{\text{F}}$  by  $\|G_1 (B_1 - U_1)\|_{\text{F}} \leq \|G_1\|_2 \|B_1 - U_1\|_{\text{F}}$  and  $\|B_1\|_{\text{F}} = \|G_1^T A_1\|_{\text{F}} \leq \|A_1\|_{\text{F}}$  by  $\|B_1\|_{\text{F}} = \|G_1^+ A_1\|_{\text{F}} \leq \|G_1^+\|_2 \|A_1\|_{\text{F}}$ .

This fact allows to introduce scaling factors in the computed carriages, e.g., in order to balance their norms, with no changes in the error estimate.

### 3.1 Backward Stability Analysis

The forthcoming Theorem 2 provides a backward stability estimate for Algorithm 3.2 which, in the exact arithmetic case, reduces to the error estimate given in Theorem 1 and, in the floating point arithmetic case, outlines the effects of the tolerance  $\varepsilon$ , the loss of orthogonality of the matrix  $M$  in (4) and the numerical stability of the function `compress` on the backward error of the computed TT decomposition.

Actually, the following analysis is mainly aimed to deal with two issues arising in the practical usage of Algorithm 3.2:

1. We compute an “exact” ( $\varepsilon = 0$ ) TT decomposition in computer arithmetics and we desire a bound on the backward error of the computed result.
2. We want to compute a “low rank” approximation ( $\varepsilon > 0$ ) of the given tensor but the function `compress` does not meet the conditions (6) as it happens if, e.g., it is based not on (pivoted) QR but on a different rank-revealing factorization where the spectral conditioning of  $M$  is greater than 1, see e.g., [1, §5.4.5]. In that case, we would like to quantify to what extent the approximation bound in Theorem 1 is degraded.

These two issues can be tackled together by assuming that the function `compress` fulfills the following hypotheses in place of (6): For any  $A \in \mathbb{R}^{m \times n}$  there exists an exact decomposition  $A + \Delta A = \widehat{Q}\widehat{R}$  with

$$\widehat{Q} = (M M') \in \mathbb{R}^{m \times m}, \quad \widehat{R} = \begin{pmatrix} N \\ N' \end{pmatrix} \in \mathbb{R}^{m \times n}, \quad (7)$$

such that  $\|M'N'\|_F \leq \varepsilon \|A\|_F$  where  $\varepsilon$  is the user-specified tolerance; and there exist two functions  $\eta_1 = \eta_1(m, n, r)$  and  $\eta_2 = \eta_2(m, r)$  such that

$$\|\Delta A\|_F \leq \eta_1 \|A\|_F, \quad \|Q - \widehat{Q}\|_2 \leq \eta_2 \quad (8)$$

for some orthogonal matrix  $Q \in \mathbb{R}^{m \times m}$ .

If computations are performed in the usual IEEE standard computer arithmetic with machine precision  $u$ , the existence of two functions  $\eta_1$  and  $\eta_2$  being  $\mathcal{O}(u)$ , fulfilling (8) and having a moderate growth in  $m, n, r$  can be inferred from known facts concerning the error analysis of Householder QR factorization [1, Ch. 5], [2, §19.3]. For example, if the matrix  $\widehat{Q}$  in (7) is computed by means of  $r \leq \min\{m, n\}$  Householder transformations then [2, pp. 359–361] brings in the estimates

$$\eta_1(m, n, r) = \mathcal{O}(mru), \quad \eta_2(m, r) = \mathcal{O}(mr^{3/2}u).$$

These estimates are deemed as rather conservative, and practical experience suggests e.g., that  $\eta_1$  is a linear polynomial in  $m$  and  $r$ , see [2, p. 368]. In any case, all quantities  $\eta_1$  and  $\eta_2$  occurring in what follows are assumed to be “sufficiently small” so that quadratic and higher order terms can be neglected. We stress that conditions (7) and (8) reduce to (6) when  $\eta_1 = \eta_2 = 0$ .



**Lemma 1.** *In the preceding notations and hypotheses, neglecting higher order terms in  $\varepsilon$ ,  $\eta_1$ , and  $\eta_2$  we have*

1.  $\|A - MN\|_F \leq (\varepsilon + \eta_1)\|A\|_F$
2.  $\|M^T(A - MN)\|_F \leq \eta_1\|A\|_F$
3.  $\|N\|_F \leq (1 + \eta_1 + \eta_2)\|A\|_F$ .

*Proof.* Firstly, we have

$$\|A - MN\|_F = \|A + \Delta A - \Delta A - MN\|_F = \|M'N' - \Delta A\|_F \leq (\varepsilon + \eta_1)\|A\|_F,$$

and the first part follows.

Let  $Q, \hat{Q} \in \mathbb{R}^{m \times m}$  be the matrices appearing in (7) and (8). Let  $Q = (Q_1 \ Q_2)$  and  $\Delta Q = \hat{Q} - Q = (\Delta_1 \ \Delta_2)$  be partitioned consistently with  $\hat{Q}$  as in (7). Then,  $\|\Delta_{1,2}\|_2 \leq \eta_2$  and

$$\begin{aligned} \|M^T M'\|_2 &= \|(Q_1 + \Delta_1)^T (Q_2 + \Delta_2)\|_2 \\ &\leq \|\Delta_1^T Q_2\|_2 + \|Q_1^T \Delta_2\|_2 + \|\Delta_1^T \Delta_2\|_2 \\ &\leq \|\Delta_1\|_2 + \|\Delta_2\|_2 + \|\Delta_1^T \Delta_2\|_2 \leq \eta_2(2 + \eta_2). \end{aligned}$$

Neglecting quadratic terms we get  $\|M^T M'\|_2 \lesssim 2\eta_2$ ,  $\|N'\|_F \lesssim \varepsilon\|A\|_F$  and

$$\begin{aligned} \|M^T(A - MN)\|_F &= \|M^T(M'N' - \Delta A)\|_F \\ &\leq \|M^T M'\|_2 \|N'\|_F + \|M\|_2 \|\Delta A\|_F \lesssim (2\eta_2\varepsilon + \eta_1)\|A\|_F. \end{aligned}$$

This proves the second inequality in the claim. Finally,

$$\begin{aligned} \|N\|_F &\leq \|\hat{R}\|_F = \|\hat{Q}^{-1}(A + \Delta A)\|_F \\ &= \|\hat{Q}^{-1}QQ^T(A + \Delta A)\|_F \\ &\leq \|\hat{Q}^{-1}(\hat{Q} - \Delta Q)\|_2 \|A + \Delta A\|_F \lesssim (1 + \eta_1)(1 + \eta_2)\|A\|_F, \end{aligned}$$

and the proof is complete.  $\square$

**Theorem 2.** *Let  $T = \text{TT}(G_1, \dots, G_d)$  where the carriages  $G_1, \dots, G_d$  are computed under the aforementioned hypotheses. Let  $\eta_1$  and  $\eta_2$  denote the largest values of the like named coefficients occurring in all recursive calls of Algorithm 3.2 on the specific input  $A$ . Moreover, let  $\hat{\varepsilon} = \varepsilon + \eta_1$ . Neglecting higher order terms in  $\varepsilon$ ,  $\eta_1$ , and  $\eta_2$ ,*

$$\|A - T\|_F \leq (1 + \eta_1 + 2\eta_2)^{d-2} ((d-2)\eta_1 + \sqrt{d-1}\hat{\varepsilon})\|A\|_F.$$

*Proof.* We proceed by induction, as in the the proof of Theorem 1. The  $d = 2$  case is an immediate consequence of Lemma 1(1):

$$\|A - T\|_F = \|A_1 - G_1 B_1\|_F \leq (\varepsilon + \eta_1)\|A_1\|_F = \hat{\varepsilon}\|A\|_F.$$

For  $d > 2$ , the inductive argument begins similarly to the proof of Theorem 1. Indeed, by hypotheses and Lemma 1 we have  $\|G_1\|_2 \leq 1 + \eta_2$  and

$$\begin{aligned} \|A - T\|_F^2 &= \|A_1 - G_1 B_1 + G_1 B_1 - G_1 U_1\|_F^2 \\ &= \|A_1 - G_1 B_1\|_F^2 + \|G_1(B_1 - U_1)\|_F^2 + 2\langle G_1^T(A_1 - G_1 B_1), B_1 - U_1 \rangle \\ &\leq \hat{\varepsilon}^2 \|A\|_F^2 + (1 + \eta_2)^2 \|B_1 - U_1\|_F^2 + 2\eta_1 \|A\|_F \|B_1 - U_1\|_F. \end{aligned}$$

The last inequality follows by Cauchy–Schwartz inequality and Lemma 1(2).

Let  $A = A_d$  and  $T = T_d$ . For  $k = 2, \dots, d-1$ , let  $A_k, T_k$  be the  $k$ -dimensional tensors defined as follows:  $A_k$  and  $T_k$  are the argument and the result of the  $(d-k)$ -th recursive call of the algorithm,<sup>1</sup>

$$(G_{d-k+1}, \dots, G_d) = \text{Full\_to\_TT}(A_k, \varepsilon), \quad T_k = \text{TT}(G_{d-k+1}, \dots, G_d).$$

For example, in the proof of Theorem 1  $A_{d-1}$  is denoted by  $B$  and  $T_{d-1}$  by  $U$ . With these notations, the preceding inequality yields

$$\|A_{k+1} - T_{k+1}\|_F^2 \leq \hat{\varepsilon}^2 \|A_{k+1}\|_F^2 + (1 + \eta_2)^2 \|A_k - T_k\|_F^2 + 2\eta_1 \|A_{k+1}\|_F \|A_k - T_k\|_F.$$

Let  $\rho = 1 + \eta_1 + \eta_2$ . From Lemma 1(3) we have  $\|A_k\|_F \leq \rho \|A_{k+1}\|_F$ . For  $k = 1, \dots, d-1$  let  $e_k = \|A_{k+1} - T_{k+1}\|_F / \|A_{k+1}\|_F$ . Neglecting product terms in  $\eta_1, \eta_2$  and other small quantities we arrive at the recurrence

$$\begin{aligned} e_k^2 &\leq \hat{\varepsilon}^2 + \rho^2 (1 + \eta_2)^2 e_{k-1}^2 + 2\eta_1 \rho e_{k-1} \\ &\leq (\alpha e_{k-1} + \eta_1)^2 + \hat{\varepsilon}^2 \end{aligned}$$

where  $\alpha = 1 + \eta_1 + 2\eta_2$  and  $e_1 \leq \hat{\varepsilon}$ . From Lemma 5 in Appendix we obtain

$$e_k \leq \alpha^{k-1} ((k-1)\eta_1 + \sqrt{k}\hat{\varepsilon}),$$

and the claim follows by setting  $k = d-1$ .  $\square$

It is worth noting that in the exact case (that is, when  $\eta_1 = \eta_2 = 0$ ) the inequality in the previous claim reduces to that of Theorem 1.

## 4 Computing Multilinear Forms

The computation of the multilinear form (or contraction)

$$\alpha = A \times_1 v^{(1)} \times_2 v^{(2)} \times_3 \dots \times_d v^{(d)}$$

---

<sup>1</sup> With a little abuse of notation, in the ensuing equation we identify  $G_{d-k+1}$  with its first unfolding matrix.

where  $v^{(i)} \in \mathbb{R}^{n_i}$ , occurs e.g., in the computation of multidimensional integrals on cartesian product grids [9]. If  $A = \text{TT}(G_1, \dots, G_d)$  then the preceding expression can be rewritten as

$$\alpha = \sum_{i_1, \dots, i_d} \sum_{j_1, \dots, j_{d-1}} G_1(i_1, j_1) G_2(j_1, i_2, j_2) \cdots G_d(j_{d-1}, i_d) v_{i_1}^{(1)} \cdots v_{i_d}^{(d)}.$$

Assuming  $n_1 = \dots = n_d = n$  and  $r_1 = \dots = r_{d-1} = r$ , the right hand side can be computed in  $\mathcal{O}(dnr^2)$  floating point operations using Algorithm 3 from [9], which is described hereafter as Algorithm 4.1.

---

**Algorithm 4.1** Fast TT contraction algorithm
 

---

**Input:** tensor train  $A = \text{TT}(G_1, \dots, G_d)$ , vectors  $v^{(1)}, \dots, v^{(d)}$   
**Output:**  $\alpha = A \times_1 v^{(1)} \times_2 \cdots \times_d v^{(d)}$   
 1: **function**  $\alpha = \text{TT\_CONTRACTION}(G_1, \dots, G_d, v^{(1)}, \dots, v^{(d)})$   
 2:      $t := G_1 \times_1 v^{(1)}$   
 3:     **for**  $k = 2 \dots d$  **do**  
 4:          $W := G_k \times_1 t$   
 5:          $t := W^T v^{(k)}$   
 6:     **end for**  
 7:      $\alpha := t$ .  
 8: **end function**

---

After completion of the  $k$ -th cycle of the for-loop,  $W$  is an  $n_k \times r_k$  matrix and  $t$  is an  $r_k$ -vector. In particular, when  $k = d$  step 4 is a matrix-vector multiplication as  $r_d = 1$ , so  $W$  becomes a vector and step 5 is an inner product of two  $n_d$ -vectors. Note that the computation of  $W = G_k \times_1 t$  followed by  $t = W^T v^{(k)}$  yields a particular algorithm to compute  $G_k \times_1 t \times_2 v^{(k)}$  which can be implemented using  $n_k r_k$  inner products with  $r_{k-1}$ -vectors followed by  $r_k$  inner products with  $n_k$ -vectors. Owing to the identity (3), an alternative algorithm for computing the same expression (with almost the same number of floating point operations) is  $W := G_k \times_2 v^{(k)}$  followed by  $t := W^T t$ .

In what follows,  $\mathbb{F}$  denotes a set of (computer) floating point numbers endowed by the usual IEEE standard arithmetics, and  $u$  denotes the corresponding unit round-off. Moreover, we use the notation  $\text{fl}(\cdot)$  with an argument that is an expression to denote the computed value of that expression. The next two lemmas are borrowed from [3].

**Lemma 2.** *Given  $x, y \in \mathbb{F}^n$ , any order of evaluation of the inner product  $x^T y$  produces an approximation  $\alpha$  such that  $|\alpha - x^T y| \leq nu|x|^T|y|$ , if no underflows or overflows are encountered.*

**Lemma 3.** *Given  $A \in \mathbb{F}^{m \times n}$  and  $x \in \mathbb{F}^n$ , let  $\hat{y} = \text{fl}(Ax)$  be the approximation to  $Ax$  obtained by computing  $m$  inner products of  $n$ -vectors, each of them being performed in an arbitrary evaluation order. If no underflow or overflow occurs, then there exists a matrix  $\hat{A} \in \mathbb{R}^{m \times n}$  such that*

$$\hat{y} = \hat{A}x, \quad \hat{A}_{ij} = (1 + \varepsilon_{ij})A_{ij}, \quad |\varepsilon_{ij}| \leq nu.$$

On the basis of these two lemmas it is not hard to obtain the result hereafter.

**Theorem 3.** *Given  $G \in \mathbb{F}^{\ell \times m \times n}$ ,  $x \in \mathbb{F}^\ell$ , and  $y \in \mathbb{F}^m$ , let  $\hat{z} = \text{fl}(\text{fl}(G \times_1 x) \times_2 y) \in \mathbb{F}^n$  be the finite precision approximation to  $z = G \times_1 x \times_2 y$  obtained after  $mn + n$  inner products, each of them being performed in an arbitrary evaluation order. If no underflow or overflow occurs then there exists a tensor  $\Delta G \in \mathbb{F}^{\ell \times m \times n}$  such that*

$$\hat{z} = (G + \Delta G) \times_1 x \times_2 y, \quad |\Delta G| \leq (\ell + m)u|G| + \mathcal{O}(u^2).$$

*Proof.* Introduce the auxiliary matrix  $M = G \times_1 x$  and let  $\hat{M} = \text{fl}(G \times_2 x)$  be its finite precision counterpart. In exact arithmetic,

$$M_{jk} = \sum_{i=1}^{\ell} G_{ijk}x_i.$$

Owing to Lemma 2, for every  $j, k$  there exists  $\varepsilon_{jk}$  such that

$$\hat{M}_{jk} - M_{jk} = \varepsilon_{jk} \sum_{i=1}^{\ell} |G_{ijk}x_i|, \quad |\varepsilon_{jk}| \leq \ell u.$$

Letting  $\eta_{ijk} = \text{sign}(G_{ijk}x_i)\varepsilon_{jk}$  we obtain

$$\begin{aligned} \hat{M}_{jk} &= \sum_i G_{ijk}x_i + \varepsilon_{jk} |G_{ijk}x_i| \\ &= \sum_i (1 + \eta_{ijk}) G_{ijk}x_i. \end{aligned}$$

Obviously  $|\eta_{ijk}| \leq \ell u$ . By Lemma 3, for  $k = 1, \dots, n$  we have

$$\begin{aligned} \hat{z}_k &= [\text{fl}(\hat{M}^T y)]_k = \sum_j (1 + \eta_{jk}) \hat{M}_{jk} y_j \\ &= \sum_j \sum_i (1 + \eta_{jk})(1 + \eta_{ijk}) G_{ijk} x_i y_j \end{aligned}$$

for some constants  $\eta_{jk}$  with  $|\eta_{jk}| \leq mu$ . In conclusion,  $\hat{z} = (G + \Delta G) \times_1 x \times_2 y$  where

$$\Delta G_{ijk} = \xi_{ijk} G_{ijk}, \quad |\xi_{ijk}| = |(1 + \eta_{jk})(1 + \eta_{ijk}) - 1| \leq (\ell + m)u + \ell mu^2,$$

and the proof is complete.  $\square$

Note that the previous theorem applies also when  $\ell = 1$  or  $n = 1$ , where  $G$  reduces to a matrix. Recalling the conventional notation  $r_0 = 1$ , we obtain immediately the following consequence.

**Corollary 1.** *Let  $\hat{\alpha}$  be the result of Algorithm 4.1 computed in machine arithmetics from input  $A = \text{TT}(G_1, \dots, G_d)$ . Then, there exist tensors  $\Delta G_1, \dots, \Delta G_d$  such that*

$$\widehat{\alpha} = \widehat{A} \times_1 v^{(1)} \times_2 v^{(2)} \cdots \times_d v^{(d)}, \quad \widehat{A} = \text{TT}(G_1 + \Delta G_1, \dots, G_d + \Delta G_d),$$

$$\text{and } |\Delta G_i| \leq (n_i + r_{i-1})u|G_i| + \mathcal{O}(u^2).$$

The previous result allows us to interpret the rounding errors in the computation of  $\alpha$  as due to perturbations in the carriages  $G_1, \dots, G_d$ , not in  $A$ . The forthcoming Theorem 4 provides a backward error bound in terms of a perturbation in  $A$ . To that goal, we need the the following lemma whose proof derives from an elementary inductive argument and will not be shown.

**Lemma 4.** *Let  $\xi_1, \dots, \xi_k$  be numbers such that  $s = \sum_{i=1}^k |\xi_i| < 1$ . Then,*

$$\prod_{i=1}^k (1 + \xi_i) = 1 + \theta, \quad |\theta| \leq \frac{s}{1-s}.$$

*In particular, if  $s \leq \frac{1}{2}$  then  $|\theta| \leq 2s$ .*

**Theorem 4.** *Let  $\widehat{\alpha}$  be the result of Algorithm 4.1 computed in machine arithmetic from input  $A = \text{TT}(G_1, \dots, G_d)$ , and let  $s = \sum_{i=1}^d (n_i + r_{i-1})u$ . If  $s \leq \frac{1}{2}$  then*

$$\widehat{\alpha} = \widehat{A} \times_1 v^{(1)} \times_2 v^{(2)} \cdots \times_d v^{(d)}$$

*where  $|\widehat{A} - A| \leq 2s \text{TT}(|G_1|, \dots, |G_d|) + \mathcal{O}(u^2)$ .*

*Proof.* In what follows, we interpret  $G_1(i, j)$  and  $G_d(i, j)$  as  $G_1(1, i, j)$  and  $G_d(i, j, 1)$ , respectively. Let  $\varepsilon_i = (n_i + r_{i-1})u$  for  $i = 1, \dots, d$ . By Corollary 1, there exist constants  $\xi_{ijk}^{(\ell)}$  such that  $\widehat{\alpha}$  is the exact contraction of the tensor  $\text{TT}(G_1 + \Delta G_1, \dots, G_d + \Delta G_d)$  and vectors  $v^{(1)}, \dots, v^{(d)}$  where

$$\Delta G_\ell(i, j, k) = G_\ell(i, j, k) \xi_{ijk}^{(\ell)}, \quad |\xi_{ijk}^{(\ell)}| \leq \varepsilon_\ell + \mathcal{O}(u^2).$$

Assuming  $j_0 = j_d = 1$ , for every multi-index  $i = (i_1, i_2, \dots, i_d)$  we have

$$\begin{aligned} \widehat{A}(i) &= \sum_{j_1, \dots, j_{d-1}} G_1(i_1, j_1) G_2(j_1, i_2, j_2) \cdots G_d(j_{d-1}, i_d) \prod_{\ell=1}^d (1 + \xi_{j_{\ell-1} i_\ell j_\ell}^{(\ell)}) \\ &= \sum_{j_1, \dots, j_{d-1}} G_1(i_1, j_1) G_2(j_1, i_2, j_2) \cdots G_d(j_{d-1}, i_d) (1 + \theta_x) \end{aligned}$$

with  $x = (i_1, \dots, i_d, j_1, \dots, j_{d-1})$  and  $|\theta_x| \leq 2s$  by Lemma 4. By triangle inequality,

$$|\widehat{A}(i) - A(i)| \leq \sum_{j_1, \dots, j_{d-1}} |G_1(i_1, j_1)| |G_2(j_1, i_2, j_2)| \cdots |G_d(j_{d-1}, i_d)| |\theta_x|,$$

and the claim follows.  $\square$

## Appendix

Hereafter, we prove a technical lemma which yields an upper bound for the growth of a sequence occurring within the proof of Theorem 2.

**Lemma 5.** *Let  $\{e_k\}$  be a sequence of nonnegative numbers such that  $e_1 \leq \gamma$  and for  $k \geq 2$*

$$e_k^2 \leq (\alpha e_{k-1} + \beta)^2 + \gamma^2$$

*for some nonnegative constants  $\alpha, \beta, \gamma$  with  $\alpha \geq 1$ . Then for all  $k \geq 1$*

$$e_k \leq \alpha^{k-1} ((k-1)\beta + \sqrt{k}\gamma).$$

*Proof.* Define the auxiliary notations  $\hat{e}_k = \alpha^{1-k} e_k$ ,  $\hat{\beta}_k = \alpha^{1-k} \beta$  and  $\hat{\gamma}_k = \alpha^{1-k} \gamma$ . Note that  $\hat{\beta}_k \leq \beta$  and  $\hat{\gamma}_k \leq \gamma$  for  $k \geq 1$ . Then,

$$\begin{aligned} \hat{e}_k^2 &= \alpha^{2-2k} e_k^2 \leq (\alpha^{2-k} e_{k-1} + \alpha^{1-k} \beta)^2 + \alpha^{2-2k} \gamma^2 \\ &= (\hat{e}_{k-1} + \hat{\beta}_k)^2 + \hat{\gamma}_k^2. \end{aligned}$$

Firstly we prove that for all  $k \geq 1$

$$\hat{e}_k \leq \sum_{j=2}^k \hat{\beta}_j + \sqrt{\sum_{j=1}^k \hat{\gamma}_j^2}.$$

Indeed, the claim is trivially verified when  $k = 1$ . By an inductive argument, for  $k \geq 2$  we have

$$\begin{aligned} \hat{e}_k^2 &\leq \left( \sum_{j=2}^k \hat{\beta}_j + \sqrt{\sum_{j=1}^{k-1} \hat{\gamma}_j^2} \right)^2 + \hat{\gamma}_k^2 \\ &= \left( \sum_{j=2}^k \hat{\beta}_j \right)^2 + \sum_{j=1}^k \hat{\gamma}_j^2 + 2 \left( \sum_{j=2}^k \hat{\beta}_j \right) \sqrt{\sum_{j=1}^{k-1} \hat{\gamma}_j^2} \\ &\leq \left( \sum_{j=2}^k \hat{\beta}_j + \sqrt{\sum_{j=1}^k \hat{\gamma}_j^2} \right)^2. \end{aligned}$$

Going back to the sequence  $\{e_k\}$  we have for all  $k \geq 1$

$$\begin{aligned} e_k &= \alpha^{k-1} \hat{e}_k \leq \alpha^{k-1} \left( \sum_{j=2}^k \hat{\beta}_j + \sqrt{\sum_{j=1}^k \hat{\gamma}_j^2} \right) \\ &\leq \alpha^{k-1} ((k-1)\beta + \sqrt{k}\gamma) \end{aligned}$$

and we are done.  $\square$

**Acknowledgements** The first author acknowledges the support received by INDAM-GNCS, Italy, for his research. The work of the second author was supported by the Russian Scientific Foundation project 14-11-00806.

## References

1. Golub, G.H., Van Loan, C.: *Matrix Computations*—4th edition. The John Hopkins University Press (2013)
2. Higham, N.J.: *Accuracy and Stability of Numerical Algorithms*. SIAM (2002)
3. Jeannerod, C.-P., Rump, S.M.: Improved error bounds for inner products in floating-point arithmetic. *SIAM J. Matrix Anal. Appl.* **34**, 338–344 (2013)
4. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM Rev.* **51**(3), 455–500 (2009)
5. Lee, N., Cichocki, A.: Estimating a few extreme singular values and vectors for large-scale matrices in tensor train format. *SIAM J. Matrix Anal. Appl.* **36**(3), 994–1014 (2015)
6. Oseledets, I.: Tensor-train decomposition. *SIAM J. Sci. Comput.* **33**(5), 2295–2317 (2011)
7. Oseledets, I., Dolgov, S.V.: Solution of linear systems and matrix inversion in the TT-format. *SIAM J. Sci. Comput.* **34**(5), A2718–A2739 (2012)
8. Oseledets, I., Tyrtshnikov, E.: Recursive decomposition of multidimensional tensors. *Dokl. Math.* **80**(1), 460–462 (2009)
9. Oseledets, I., Tyrtshnikov, E.: TT-cross approximation for multidimensional arrays. *Lin. Alg. Appl.* **432**(1), 70–88 (2010)