

# From low rank tensors for grids of multivariate functions to joint distributions of statistical models: selecting the relevant norm.

Alain Franc

INRAE BioGeCo & INRIA Pleiade

July 8, 2021

INRIA Skoltech workshop



# Three objects behind a tensor

$$\mathbf{A} \in E_1 \otimes \dots \otimes E_d \equiv \mathbb{R}^{n_1 \times \dots \times n_d}$$

discretization  
of a multivariate  
function

joint law of  
a statistical  
model

array  
of data

# Preliminaries

## CP decomposition

$$\mathbf{A} = \sum_{\alpha=1}^r \mathbf{x}_{\alpha} \otimes \mathbf{y}_{\alpha} \otimes \mathbf{z}_{\alpha} \quad \in \mathbb{R}^{n \times n \times n}$$

$$r_{\text{cp}}(\mathbf{A}) = \inf \left\{ r \in \mathbb{N} \mid \exists \mathbf{A} = \sum_{\alpha=1}^r \dots \right\}$$

$3nr$  terms

## TT decomposition

$$\mathbf{A}[i, j, k] = \mathbf{u}[i] \cdot \mathbf{G}[j] \cdot \mathbf{v}[k]$$

$$r_{\text{tt}} : \mathbf{u}_i \in \mathbb{R}^{1 \times r}, \quad \mathbf{G}_j \in \mathbb{R}^{r \times r}, \quad \mathbf{v}_k \in \mathbb{R}^{r \times 1}$$

$nr^2 + 2nr$  terms

Lemma (can be generalized in a straightforward way)

Let

$$\mathbf{A} \in E \otimes F \otimes G, \quad \text{with} \quad \begin{cases} \dim E & = m \\ \dim F & = n \\ \dim G & = p \end{cases}$$

Then

$$r_{tt}(\mathbf{A}) \leq r_{cp}(\mathbf{A}) \leq r_{tt}^2(\mathbf{A})$$

Proof

- Development on a basis and reorganisation for  $r_{cp} \leq r_{tt}^2$
- Development with diagonal  $G$  matrix for  $r_{tt} \leq r_{cp}$

# Example

For example, if  $r = 2$

$$\mathbf{A} = \mathbf{x}_1 \otimes \mathbf{y}_1 \otimes \mathbf{z}_1 + \mathbf{x}_2 \otimes \mathbf{y}_2 \otimes \mathbf{z}_2$$

$$\text{with } \mathbf{x}_1 = (x_i^{(1)})_i, \quad \mathbf{y}_1 = (y_j^{(1)})_j, \dots$$

Then

$$\begin{aligned} x_i^{(1)} y_j^{(1)} z_k^{(1)} + x_i^{(2)} y_j^{(2)} z_k^{(2)} &= \begin{pmatrix} x_i^{(1)} & x_i^{(2)} \end{pmatrix} \begin{pmatrix} y_j^{(1)} & z_k^{(1)} \\ y_j^{(2)} & z_k^{(2)} \end{pmatrix} \\ &= \begin{pmatrix} x_i^{(1)} & x_i^{(2)} \end{pmatrix} \begin{pmatrix} y_j^{(1)} & 0 \\ 0 & y_j^{(2)} \end{pmatrix} \begin{pmatrix} z_k^{(1)} \\ z_k^{(2)} \end{pmatrix} \\ &= \mathbf{u}_i \cdot \mathbf{G}_j \cdot \mathbf{v}_k \end{aligned}$$

When  $r_{cp} = r_{tt}$ ?

TT rank is a tensor property (whatever the basis), and

$$\begin{aligned}u_i \cdot G_j \cdot v_k &= u_i \cdot (PP^{-1}) \cdot G_j \cdot (QQ^{-1}) \cdot v_k \\&= (u_i \cdot P)(P^{-1} \cdot G_j \cdot Q)(Q^{-1} \cdot v_k) \\&= u'_i G'_j v'_k\end{aligned}$$

Simultaneously diagonalizable matrices

- $(G_j)_{1 \leq j \leq n}$  is a set of  $r \times r$  **diagonalizable** matrices

Then, it is equivalent that

- $\forall (i, j), \quad G_i G_j = G_j G_i$
- There exists a basis  $P$  such that each  $G_j$  is diagonal in this basis

## Lemma (refined)

Let  $\mathbf{A} \in \mathbb{R}^{m \times n \times p}$  with  $a_{ijk} = u_i \cdot G_j \cdot v_k$ ,

### Lemma

- $r_{\text{tt}} \leq r_{\text{cp}} \leq r_{\text{tt}}^2$
- $r_{\text{cp}} = r_{\text{tt}}$  iff  $(G_j)_j$  simultaneously diagonalisable

### Take home message: very simple

If a tensor has low CP rank, then it has low TT rank.



# A tensor as discretisation of a multivariate function

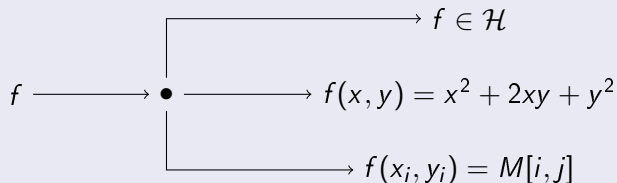
# Three objects behind a function

Multivariate function (here  $d=2$  for simplicity)

Let

$$\mathbb{R}^2 \xrightarrow{f} \mathbb{R}$$

Three implementations



## 1D mesh

$$\mathbf{x} \in \mathbb{R}^n, \quad \forall i, \quad x_i \in \mathbb{R}$$

$$\mathbf{x} = (x_1, \dots, x_n) \quad \text{with} \quad x_1 < \dots < x_n$$

$$\mathcal{D}(f, \mathbf{x}) = (f(x_1), \dots, f(x_n))$$

$$\mathcal{H} \times \mathcal{M} \xrightarrow{\mathcal{D}} \mathbb{R}^n$$

## 2D Cartesian Mesh

- Let us consider 2 meshes  $\mathbf{x}, \mathbf{y}$  in  $\mathbb{R}$  of respective sizes  $m, n$ .
- Then  $\mathbf{x} \otimes_m \mathbf{y}$  in  $\mathbb{R}^2$  can be defined as

$$\mathbf{x} \otimes_m \mathbf{y} = \begin{pmatrix} (x_1, y_1) & (x_1, y_2) & \dots & (x_1, y_n) \\ (x_2, y_1) & (x_2, y_2) & \dots & (x_2, y_n) \\ \vdots & \vdots & & \vdots \\ (x_m, y_1) & (x_m, y_2) & \dots & (x_m, y_n) \end{pmatrix}$$

# Discretization of product of functions

## Product of functions: two variables

- $(fg)(x) = f(x)g(x)$
- $(f \otimes g)(x, y) = f(x)g(y)$

## Discretisation of product of functions

- $\mathcal{D}(fg, \mathbf{x}) = \mathcal{D}(f, \mathbf{x}) \odot \mathcal{D}(g, \mathbf{x})$
- $\mathcal{D}(f \otimes g, \mathbf{x} \otimes_m \mathbf{y}) = \mathcal{D}(f, \mathbf{x}) \otimes \mathcal{D}(g, \mathbf{y})$

## Consequence

If  $u = f \otimes g$  [ $u(x, y) = f(x)g(y)$ , separation of variables], the matrix  $\mathcal{D}(u, \mathbf{x} \otimes_m \mathbf{y})$  of discretization of  $u$  on mesh  $\mathbf{x} \otimes_m \mathbf{y}$  has rank one too.

Generalization is straightforward ...

$$\mathcal{D} \left( \prod_{\mu} f_{\mu}, \mathbf{x} \right) = \odot \mathcal{D}(f_{\mu}, \mathbf{x})$$

$$\mathcal{D} \left( \bigotimes_{\mu} f_{\mu}, \bigotimes_{\mu}^{(m)} \mathbf{x}_{\mu} \right) = \bigotimes_{\mu} \mathcal{D}(f_{\mu}, \mathbf{x}_{\mu})$$

# An elementary lemma (can be extended to $d > 2$ )

Let

$$\begin{aligned}\mathcal{H} \otimes \mathcal{H} &\xrightarrow{\mathcal{D}} \mathbb{R}^{m \times n} \\ \psi &\longrightarrow \mathcal{D}(\psi, \mathbf{x} \otimes_m \mathbf{y})\end{aligned}$$

Let (rank  $r$  CP-decomposition)

$$\psi = \sum_{\alpha=1}^r \mathbf{u}_\alpha \otimes \mathbf{v}_\alpha \quad \text{with} \quad \mathbf{u}_\alpha, \mathbf{v}_\alpha \in \mathcal{H}$$

Then

$$\forall m, n, \quad r_{\text{cp}}(\mathcal{D}(\psi, \mathbf{x} \otimes_m \mathbf{y})) \leq r_{\text{cp}}(\psi)$$

As consequence ...

$$\forall m, n, \quad r_{\text{tt}}(\mathcal{D}(\psi, \mathbf{x} \otimes_m \mathbf{y})) \leq r_{\text{cp}}(\psi)$$

# Discretization of polynomials (1/2)

Let us define

$$\mathbb{R} \xrightarrow{\mathbf{1}_x} \mathbb{R}$$

$$x \longrightarrow 1$$

Bivariate polynomial

$$P(x, y) = x^2 + 2xy + y^2$$

Then

$$P = x^2 \otimes \mathbf{1}_y + 2x \otimes y + \mathbf{1}_x \otimes y^2$$

and

$$r_{\text{cp}}(P) \leq 3$$

and

$$\forall m, n \in \mathbb{N}, \quad r_{\text{tt}}(\mathcal{D}(P, \mathbf{x} \otimes_m \mathbf{y})) \leq r_{\text{cp}}(\mathcal{D}(P, \mathbf{x} \otimes_m \mathbf{y})) \leq 3$$

# Discretization of polynomials (2/2)

## Any polynomial

$$\mathbf{P}(x_1, \dots, x_d) = \sum_{\mathbf{n}} a_{n_1 \dots n_d} x_1^{n_1} \dots x_d^{n_d}, \quad N \text{ terms}$$

Then

$$\mathbf{P} = \sum_{\mathbf{n}} a_{n_1 \dots n_d} x_1^{n_1} \otimes \dots \otimes x_d^{n_d}, \quad \text{CP rank} = N$$

and

$$\forall m_1, \dots, m_d \in \mathbb{N}, \quad r_{\text{tt}}(\mathcal{D}(\mathbf{P}, \mathbf{x}_1 \otimes_m \dots \otimes_m \mathbf{x}_d)) \leq N$$

- $m_\mu$  is the size of the Cartesian Grid for mode  $\mu$ .

## a remark

In general,  $N = \prod_{\mu} n_{\mu}$

$\implies$  rank is high if  $d$  is significant



## Stone-Weierstrass theorem

- $C(S) = \{f : S \rightarrow \mathbb{R}, f \text{ continuous}\}$
- $A \subset C(S)$  s.t.  $f, g \in A \Rightarrow fg \in A$
- $x \neq y \Rightarrow \exists f \in A : f(x) \neq f(y)$

Then

$$\forall \epsilon > 0, \forall f \in C(S), \exists \varphi \in A : \forall x \in S, |f(x) - \varphi(x)| < \epsilon$$

- Any continuous function can be approximated with norm  $\ell^\infty$  as close as wished by a polynomial, i.e. a low rank function
- This is automatically transported to approximation by low rank discretization on Cartesian grids
- And from CP to TT approximation (for working with TT toolbox)
- A well developed theory has been elaborated for  $\ell^2$  norm as well (development on basis of **orthogonal polynomials**, leading to **Tucker approximations**).

# Take home message

There is a sound and standard algebraic theory for showing that tensors as discretization of multivariate functions are

- exactly low rank for polynomials
- well approximated by low rank tensors (CP, TT, Tucker) for continuous functions

## Some details still deserve attention

- Better understanding of the link between CP et TT decomposition
  - Why CP is numerically unstable sometimes and TT not ?
  - at which boundaries ?
- is it expandable to non Cartesian meshes?

# A tensor as a joint law

# A tensor as a joint law

## Setting

- Let us have a discrete set  $\Lambda$
- $d$  random variables  $X_\mu$  with values in  $\Lambda$

Define

$$\mathbf{T}[i_1, \dots, i_d] = t_{i_1 \dots i_d} \propto \mathbb{P}(X_1 = i_1, \dots, X_d = i_d)$$

- $\mathbf{T}$  is a  $d$ -modes tensor with elements  $\geq 0$

## Classical examples

- Ising model
- Graphical models
- ...

## Definition

In statistical physics, and statistical modeling, one is led to compute

$$Z(\mathbf{T}) = \sum_{i_1} \dots \sum_{i_d} t_{i_1 \dots i_d}$$

- requires  $n^d$  additions, with  $d > 10^3$  often .... see Novikov & al. (2014) for large  $d$ .

Alexander Novikov & al. (2014) - Putting MRFs on a Tensor Train.  
*Proceedings of the 31 st International Conference on Machine Learning*,  
Beijing, China.

# When does it work?

## When $\mathbf{T}$ is written in TT format

$$t_{i_1 \dots i_d} = G_1(i_1) \times \dots \times G_\mu(i_\mu) \times \dots \times G_d(i_d)$$

with  $G_\mu(i_\mu) \in \mathbb{R}^{r_{\mu-1} \times r_\mu}$ . Then, if  $B_\mu = \sum_{i_\mu} G_\mu(i_\mu)$

$$Z(\mathbf{T}) = B_1 B_2 \dots B_d$$

(Novikov & al., 2014)

## When $\mathbf{T}$ has low CP-rank

$$\mathbf{T} = \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{z} + \mathbf{x}' \otimes \mathbf{y}' \otimes \mathbf{z}'$$

## Observation

- All terms in  $\mathbf{T}$  are  $\geq 0$  Then

$$Z(\mathbf{T}) = \|\mathbf{T}\|_1 \quad (\text{norm } \ell^1)$$

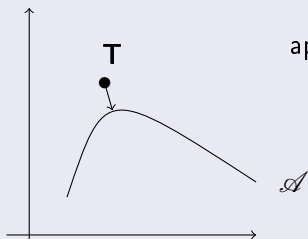
## General approach

- a quantity  $f(\mathbf{T})$  on a tensor (like  $Z(\mathbf{T})$ ) is difficult to compute in general, but easy in a (closed) subset  $\mathcal{A} \subset \mathbb{R}^{n_1 \times \dots \times n_d}$
- Then, one computes  $\hat{\mathbf{T}} \in \mathcal{A}$  such that, for a selected distance  $\delta$

$$\delta(\mathbf{T}, \hat{\mathbf{T}}) \text{ is minimal}$$

- and approximates

$$f(\mathbf{T}) \approx f(\hat{\mathbf{T}})$$



approaching  $Z(\mathbf{T})$  with  $\ell^1$  norm

## Setting the problem

- given  $A \in \mathbb{R}^{m \times n}$
- find  $\hat{A}$  with  $\text{rank}(\hat{A}) = r$
- such that  $\|A - \hat{A}\|_1$  minimal

## State of the art

- Difficulty: it has been shown to be NP-hard (Gillis & Vavasis, 2015)
- Available algorithm with "provable approximation guarantees" (Song, Woodruff & Zhong, 2018)



# What about tensors?

- Let  $\mathbf{A} \in E \otimes F \otimes G$
- Let  $A_E$  be its first matricization

$$F \otimes G \xrightarrow{A_E} E$$

- Then  $\|A_E\|_1 = \|\mathbf{A}\|_1$

## Proposed heuristics

for all matricizations  $A$  of  $\mathbf{A}$  do

**computes** the best rank one approximation  $A_{best}$  of  $A$  with norm  $\ell^1$

**computes**  $\delta(A) = \|A - A_{best}\|$

select  $A$  such that  $\delta(A)$  is minimal

**computes**  $Z(A_{best})$  (easy)

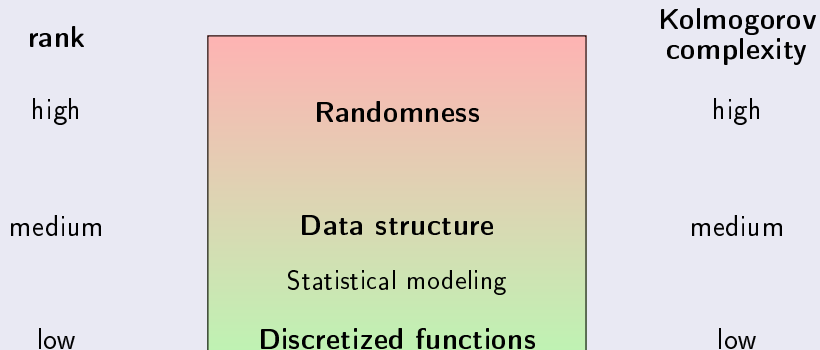
**return**  $Z(A_{best}) \approx Z(\mathbf{A})$

# Take home message

- The theory behind best low rank approximation of a tensor as joint distribution is not as mature than links between multilinear algebra, PDE, functional analysis
- A difficulty is that a wider diversity of norms is relevant
  - Kullback-Leibler for mutual information
  - $\ell^1$  for partition function
  - ...

which do not rely on Euclidean geometry

# To be taken into account: different levels of complexity



# Acknowledgements

Thanks to following colleagues for helpful discussions

- Mohamed Anwar Abouabdallah, PhD Student
- Olivier Coulaud, INRIA Hiepacs
- Martina Iannacito, PhD student
- Nathalie Peyrard, INRAE MIAT, Toulouse

and Olivier Beaumont for the proposition to join the team!