# Low-rank approximation: from matrices to tensors

Konstantin Usevich, CNRS, CRAN (UMR 7039 / Univ. Lorraine ), Nancy

31.05.2024, Marseille, workshop on Low-rank optimization



#### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

Extensions

## Overview

#### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

Extensions

## Matrix rank: factorization view

Rank of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  (or  $\mathbb{C}^{m \times n}$ )

 $\blacktriangleright$  = number of linearly independent columns/ rows in  ${\bf X}$ 

# Matrix rank: factorization view

Rank of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  (or  $\mathbb{C}^{m \times n}$ )

Matrix LRA

- $\blacktriangleright \stackrel{\rm def}{=}$  number of linearly independent columns/ rows in  ${\bf X}$
- $\blacktriangleright$  = minimal r such that X can be factorized as



 $\Rightarrow$  rank is bounded by dimensions:

 $\operatorname{\mathsf{rank}} \mathbf{X} \le \min(m, n)$ 



Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ):





Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ):



Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ):

Matrix LRA



Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ):

Matrix LRA



Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ):

Matrix LRA





Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ), cont'd:





Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ), cont'd:



(Scotland)

(Wales)



(Scotland) rank  $\approx \frac{m}{2}$  (symmetry)

Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m imes n}$ ), cont'd:



(Wales)



Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m imes n}$ ), cont'd:



(Scotland) rank  $pprox rac{m}{2}$  (symmetry)



(Wales) 
$${\sf rank} pprox rac{5}{6}m$$
 (finite support)

Ranks of flags (as  $\mathbf{X} \in \mathbb{R}^{m imes n}$ ), cont'd:



Matrix LRA

$$({\sf Scotland})$$
 rank  $pprox rac{m}{2}$  (symmetry)

(Wales) rank  $\approx \frac{5}{6}m$  (finite support)

- random matrix (with a.c. probability distribution): rank = min(m, n) a.s. (with probability 1)
- many interesting matrices are [well approximated by] low-rank [Townsend, Udell, 2017]

# Singular value decomposition (SVD)

a/the ("economy size") SVD ( $m \le n$ ):

$$m \boxed{\boldsymbol{X}}^{n} = \begin{bmatrix} \boldsymbol{U} \end{bmatrix}^{\sigma_{1}} \vdots \vdots \\ \sigma_{m} \begin{bmatrix} \boldsymbol{V}^{T} \end{bmatrix} = \sum_{k=1}^{m} \sigma_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{T}$$

where

Matrix LRA

- $U^T U = V^T V = I$  semi-orthogonal matrices of singular vectors  $U = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_m \end{bmatrix}$ ,  $V = \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_m \end{bmatrix}$
- $\sigma_1 \geq \cdots \geq \sigma_m \geq 0$  singular values

# Singular value decomposition (SVD)

a/the ("economy size") SVD ( $m \le n$ ):

$$m \boxed{\boldsymbol{X}}^{n} = \begin{bmatrix} \boldsymbol{U} \end{bmatrix}^{\sigma_{1}} \vdots_{\sigma_{m}} \begin{bmatrix} \boldsymbol{V}^{T} \end{bmatrix} = \sum_{k=1}^{m} \sigma_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{T}$$

where

Matrix LRA

•  $U^T U = V^T V = I$  — semi-orthogonal matrices of singular vectors  $U = \begin{bmatrix} u_1 & \cdots & u_m \end{bmatrix}, \quad V = \begin{bmatrix} v_1 & \cdots & v_m \end{bmatrix}$ 

•  $\sigma_1 \geq \cdots \geq \sigma_m \geq 0$  — singular values



# SVD and low-rank approximation

#### Eckart-Young(-Mirsky(-Schmidt)) theorem: best rank-r approximation

$$\min_{\widehat{\mathbf{X}}} \| \widehat{\mathbf{X}} - \mathbf{X} \|$$
 subject to rank  $\widehat{\mathbf{X}} \leq r$ 

in any unitarily invariant\* norm  $\|\cdot\|$  is given by

$$\mathsf{tSVD}_r(\boldsymbol{X}) :=$$
 $\boldsymbol{U} \begin{bmatrix} \sigma_1 \\ \vdots \\ \vdots \\ \sigma_r \end{bmatrix} \quad \boldsymbol{V}^T$  $\boldsymbol{0}$ 

Matrix LRA

truncated SVD Singular Values  
Hard  
Thresholding  
\* Examples: Frobenius norm 
$$\|\mathbf{X}\|_F^2 := \sum_{i,j=1}^{m,n} X_{ij}^2$$
, spectral norm...

# SVD and low-rank approximation

#### Eckart-Young(-Mirsky(-Schmidt)) theorem: best rank-r approximation

$$\min_{\widehat{\mathbf{X}}} \| \widehat{\mathbf{X}} - \mathbf{X} \| \quad \text{subject to} \quad \text{rank } \widehat{\mathbf{X}} \leq r$$

in any unitarily invariant\* norm  $\|\cdot\|$  is given by

Matrix LRA

$$\mathsf{tSVD}_{r}(\boldsymbol{X}) := \boxed{\boldsymbol{U}} \stackrel{\sigma_{1}}{\stackrel{\cdots}{\overset{\sigma_{1}}{\odot}}}_{\boldsymbol{0}} \boxed{\boldsymbol{V}^{T}} = \boxed{\boldsymbol{U}_{1:r,:}} \stackrel{\sigma_{1}}{\stackrel{\cdots}{\overset{\sigma_{1}}{\odot}}} \boxed{\boldsymbol{V}_{1:r,:}^{T}} = \sum_{k=1}^{r} \sigma_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{T}$$

$$\mathsf{truncated SVD} \stackrel{\mathsf{Singular Values}}{\stackrel{\mathsf{Hard}}{\overset{\mathsf{Thresholding}}{\overset{\mathsf{Hard}}{\overset{\mathsf{Thresholding}}{\overset{\mathsf{Thresho$$

Low-rank approximations: example







Matrix LRA

100

200



Computational aspects: full and partial SVD

$$m \boxed{\boldsymbol{X}}^{n} = \begin{bmatrix} \boldsymbol{U} & \sigma_{1} \\ \ddots & \sigma_{m} \end{bmatrix} \boldsymbol{V}^{T} = \sum_{k=1}^{m} \sigma_{k} \mathbf{u}_{k} \mathbf{v}_{k}^{T}$$

Full SVD  $(m \le n)$ :  $\mathcal{O}(m^2 n)$  ([Golub, Reinsh, 1970])

•  $(\mathbf{u}_k, \sigma_k^2)$  — eigenvalues of  $\mathbf{C} = \mathbf{X}\mathbf{X}^T$ :  $\mathcal{O}(m^3)$  if  $m \ll n$ 

Truncated SVD (r first eigenvalues/vectors):

Matrix LRA

- Iterative (e.g., Lanzos/Arnoldi) algorithms (e.g., [Simon, 1984]) "generalization of power method"
- Randomized SVD [Halko, Martinsson, Tropp, 2011]

Both roughly |O(Mr)|, where  $M \leq mn$  (see next slide):

## Power iteration

Goal: find the top eigenpair  $\mathbf{u}_1, \lambda_1$  of  $\mathbf{C}$  .

**>** Set 
$$\mathbf{u}^{(0)} \in \mathbb{R}^{m \times m}$$
 random.

• Iterate 
$$\mathbf{u}^{(k+1)} = \frac{\mathbf{C}\mathbf{u}^{(k)}}{\|\mathbf{C}\mathbf{u}^{(k)}\|}$$

Matrix LRA

Case  $\mathbf{C} = \mathbf{X}\mathbf{X}^T$ : cost  $\approx \mathcal{O}(Mn_{iter})$ , where

• full matrices: 
$$M = mn$$

- ▶ sparse matrices: M = #non-zero entries
- structured matrices: (e.g.,  $M = n \log(n)$  for Hankel)

M = cost of matrix-vector product (e.g.,  $\mathbf{X}\mathbf{v}$ )

- used e.g., by Google for PageRank
- do not need to store the matrix X
- generalizes to rank-r approximation (cost  $\mathcal{O}(Mr + mr^2)$ )

**Observation.** A rank-r matrix is uniquely determined by the  $r \times r$  cross if the  $r \times r$  submatrix  $\mathbf{X}_{\mathcal{I},\mathcal{J}}$  is nonsingular

$$\mathbf{X} = \mathbf{X}_{:,\mathcal{J}} (\mathbf{X}_{\mathcal{I},\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{I},:}$$
(\*)

**Example.** r = 2:

Matrix LRA



**Observation.** A rank-r matrix is uniquely determined by the  $r \times r$  cross if the  $r \times r$  submatrix  $\mathbf{X}_{\mathcal{I},\mathcal{J}}$  is nonsingular

$$\mathbf{X} = \mathbf{X}_{:,\mathcal{J}} (\mathbf{X}_{\mathcal{I},\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{I},:}$$
(\*)

**Example.** r = 2:

Matrix LRA



**Observation.** A rank-r matrix is uniquely determined by the  $r \times r$  cross if the  $r \times r$  submatrix  $\mathbf{X}_{\mathcal{I},\mathcal{J}}$  is nonsingular

$$\mathbf{X} = \mathbf{X}_{:,\mathcal{J}} (\mathbf{X}_{\mathcal{I},\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{I},:}$$
(\*)

**Example.** r = 2:

Matrix LRA



**Observation.** A rank-r matrix is uniquely determined by the  $r \times r$  cross if the  $r \times r$  submatrix  $\mathbf{X}_{\mathcal{I},\mathcal{J}}$  is nonsingular

$$\mathbf{X} = \mathbf{X}_{:,\mathcal{J}} (\mathbf{X}_{\mathcal{I},\mathcal{J}})^{-1} \mathbf{X}_{\mathcal{I},:}$$
(\*)

**Example.** r = 2:

Matrix LRA



 $\Rightarrow$  (\*) can be used as an approximation

If the matrix is huge or expensive to compute:

CUR (cross, or pseudo-skeleton) approximation (for size-r subsets  $\mathcal{I}, \mathcal{J}$ ):

$$\widehat{\mathbf{X}}_{cross}(\mathcal{I},\mathcal{J}) = \mathbf{X}_{:,\mathcal{J}}(\mathbf{X}_{\mathcal{I},\mathcal{J}})^{-1}\mathbf{X}_{\mathcal{I},:}$$

[Mahoney, Drineas, 2012] Advantages:

Matrix LRA

- ▶ Need a small portion (cross) of the matrix ( O(r(m+n)))
- Quasi-optimality (thm. in [Goreinov, Tyrtyshnikov, 2001])

$$\|\widehat{\mathbf{X}}_{cross}^* - \mathbf{X}\|_{max} \le (r+1)\sigma_r(\mathbf{X})$$

for  $|\det(\mathbf{X}_{\mathcal{I},\mathcal{J}})| \to \max$ 

• iterative or randomized strategies to select  $\mathcal{I}, \mathcal{J}$ 

## Extensions of the basic problem

$$\min_{\substack{\widehat{\mathbf{X}} = \mathbf{AB} \\ \mathbf{A} \in \mathbb{R}^{m \times r}, \mathbf{B} \in \mathbb{R}^{r \times n}}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_{F}$$

• Other norms (e.g. 
$$\|\mathbf{X}\|_W^2 = \sum_{i,j} W_{ij} X_{ij}^2$$
), missing data

 Constraint on the matrix: structured X — structured low-rank approximation

Matrix LRA

Constraints on the factors A, B (e.g., nonnegative)

# Weighted (unstructured) LRA

$$\min_{\widehat{\mathbf{X}}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_W^2 \quad \text{subject to} \quad \operatorname{rank}(\widehat{\mathbf{X}}) \leq r$$

where  $\|\mathbf{X}\|_W^2 = \sum_{i,j=1}^{m,n} \mathbf{X}_{ij}^2 W_{ij}, \quad W_{ij} \in (0; +\infty)$  weighted norm

• 
$$W_{ij} \equiv 1$$
 (or rank  $W = 1$ )  $\rightarrow$  solution by SVD

Matrix LRA

0000000000000000

In general case, no closed form solution: Gillis, Glineur, Low-Rank Matrix Approximation with Weights or Missing Data is NP-hard, SIMAX, 2011.

## Extended semi-norms and matrix completion

$$\label{eq:constraint} \min_{\widehat{\mathbf{X}}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_W^2 \text{ subject to } \quad \mathrm{rank}(\widehat{\mathbf{X}}) \leq r$$

 $\|\mathbf{X}\|_W^2 := \sum_{i,j=1}^{m,n} \mathbf{X}_{ij}^2 W_{ij}, \ W_{ij} \in [0; +\infty]$  weighted extended semi-norm

• fixed values: 
$$W_{i,} = +\infty \longleftrightarrow$$
 constraint  $\mathbf{X}_{ij} = \widehat{\mathbf{X}}_{ij}$ 

• missing values:  $W_{ij} = 0 \longleftrightarrow \widehat{\mathbf{X}}_{ij}$  is not important

Matrix LRA

## Extended semi-norms and matrix completion

$$\min_{\widehat{\mathbf{X}}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_W^2 \text{ subject to } \quad \operatorname{rank}(\widehat{\mathbf{X}}) \leq r$$

 $\|\mathbf{X}\|_W^2 := \sum_{i,j=1}^{m,n} \mathbf{X}_{ij}^2 W_{ij}, \ W_{ij} \in [0;+\infty] \quad \text{ weighted extended semi-norm}$ 

• fixed values: 
$$W_{i,} = +\infty \longleftrightarrow$$
 constraint  $\mathbf{X}_{ij} = \widehat{\mathbf{X}}_{ij}$ 

• missing values:  $W_{ij} = 0 \longleftrightarrow \widehat{\mathbf{X}}_{ij}$  is not important

Extreme case:  $W_{ij} \in \{0, +\infty\}$  — exact matrix completion

 $\min \operatorname{rank}(\widehat{\mathbf{X}})$ subject to  $(\widehat{X})_{ij} = X_{ij}, \quad \forall (i,j) \in \Omega$ 

Matrix LRA



# Structured low-rank approximation

Matrix LRA

[Markovsky, 2008] : Problem (SLRA). Given a structured matrix  $\mathbf{X} \in \mathcal{S}$ 

$$\underset{\widehat{X}}{\mathsf{minimize}} ~~ \|\mathbf{X} - \widehat{\mathbf{X}}\|_W^2 ~~ \mathsf{subject to} ~~ \widehat{\mathbf{X}} \in \mathcal{S} ~~ \mathsf{and} ~~ \mathsf{rank} ~ \widehat{\mathbf{X}} \leq r$$

 $\mathsf{Data}\approx\mathsf{low-complexity}\ \mathsf{model}$ 

structure ${\cal S}$	approximation problem
unstructured	fit by <i>r</i> -dim. subspace
Hankel	fitting by complex exponentials
block-Hankel	linear system identification
	model reduction
Sylvester	approx. greatest common divisor
generalized	fit set of points by
Vandermonde	algebraic hypersurfaces

Source separation  $\leftrightarrow$  matrix factorization

Example from spectroscopy:

Matrix LRA

00000000000000

- each observed spectrum is a linear combination of "pure spectra"
- different conditions different coefficients.



Instantaneous mixture model:  $\mathbf{X}(x,\lambda) = \sum_k \mathbf{a}_k(x) \mathbf{s}_k(\lambda)$ 

Source separation  $\leftrightarrow$  matrix factorization

Example from spectroscopy:

000000000000

Matrix LRA

- each observed spectrum is a linear combination of "pure spectra"
- different conditions different coefficients.



Instantaneous mixture model:  $\mathbf{X}(x,\lambda) = \sum_k \mathbf{a}_k(x) \mathbf{s}_k(\lambda)$ 



Can we recover  $\mathbf{A}$ ,  $\mathbf{S}$  from  $\mathbf{X}$ ?

Matrix factorizations are non-unique

Does not happen for matrices:

Matrix LRA

$$\mathbf{M} = \mathbf{A} \mathbf{S}^{T} = \mathbf{A} \mathbf{Q} \mathbf{Q}^{T} \mathbf{S}^{T}$$

nonunique (change of basis)

However, constraints on factors can guarantee essential uniqueness

$$\mathbf{Q} = \mathbf{\Lambda} \cdot \mathbf{\Pi}$$

diagonal permutation

Examples of constraints:





# Other tools/link to optimization

 $\blacktriangleright \ \mathcal{M}_r^{m \times n} = \{ \mathbf{X} \in \mathbb{R}^{m \times n} | \operatorname{rank}(\mathbf{X}) = r \} - \text{smooth manifold}$ 



visualisation of SLRA

 $\rightarrow$  optimization on manifolds [Absil, Mahoney, Sepulchre, 2008], [Boumal, 2023]

- ▶  $\mathcal{M}_{\leq r}^{m \times n} = { \mathbf{X} | \operatorname{rank}(\mathbf{X}) \leq r }$  algebraic variety (stratified set) link to determinantal representations of algebraic varieties
- $\blacktriangleright \text{ low rank} \leftrightarrow \text{sparsity of singular values}$

Matrix LRA

000000000000

$$\underbrace{\|(\sigma_1,\ldots,\sigma_m)\|_0}_{\text{rank}}\leftrightarrow\underbrace{\|(\sigma_1,\ldots,\sigma_m)\|_1}_{\text{nuclear norm}}$$

other sparsity-promoting penalties ...

- other norms/divergences/losses....
- dynamical low-rank approximation  $\mathbf{X}(t) \in \mathcal{M}_r$ , varies over time
### Matrix factorization $\leftrightarrow$ neural networks

Matrix LRA

0000000000000000 00000

Matrix factorization:  $\mathbf{X} = \mathbf{U}\mathbf{V}^T, \mathbf{U} = \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_r \end{bmatrix}, \mathbf{V} = \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_r \end{bmatrix}$ Decompose linear map  $\mathbf{f} : \mathbb{R}^m \to \mathbb{R}^n, \mathbf{f}(\mathbf{z}) = \mathbf{X}\mathbf{z}$  as

 $\mathbf{f}(\mathbf{z}) = \mathbf{u}_1 \cdot (\mathbf{v}_1^\top \mathbf{z}) + \dots + \mathbf{u}_r \cdot (\mathbf{v}_r^\top \mathbf{z}),$ 

### Matrix factorization $\leftrightarrow$ neural networks

Given nonlinear map  $\mathbf{f}:\mathbb{R}^m \to \mathbb{R}^n$ , decompose it as

Matrix LRA

0000000000000000

$$\mathbf{f}(\mathbf{z}) = \mathbf{u}_1 \underline{g}_1(\mathbf{v}_1^\top \mathbf{z}) + \dots + \mathbf{u}_r \underline{g}_r(\mathbf{v}_r^\top \mathbf{z}),$$

where  $g_k(t)$  are univariate functions (see. e.g., [Comon,Qi,U., 2017])

### Matrix factorization $\leftrightarrow$ neural networks

Given nonlinear map  $\mathbf{f}:\mathbb{R}^m\to\mathbb{R}^n,$  decompose it as

$$\mathbf{f}(\mathbf{z}) = \mathbf{u}_1 \underline{g}_1(\mathbf{v}_1^\top \mathbf{z}) + \dots + \mathbf{u}_r \underline{g}_r(\mathbf{v}_r^\top \mathbf{z}),$$

where  $g_k(t)$  are univariate functions (see. e.g., [Comon,Qi,U., 2017])



See for example:

Matrix LRA

0000000000000000

- One-hidden layer model: ([Marcotte, Gribonval, Peyré, 2024])
- deep linear networks (products of matrices): [Malgouyres, 2020]
- deep NMF [Leplat et al, 2024]

# Overview

### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

Extensions

### Tensors: some notation

 $n_{d}$ 

Tensor product of vector spaces (over a field  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{C}$  ):

▶ 
$$\mathcal{T} \in \mathbb{F}^{n_1 \times \dots \times n_d} \stackrel{\text{def}}{=} \mathbb{F}^{n_1} \otimes \mathbb{F}^{n_2} \otimes \dots \otimes \mathbb{F}^{n_d}$$
  
▶ *d*-way array  $\mathcal{T} = [\mathcal{T}_{i_1, i_2, \dots, i_d}]_{i_1, i_2, \dots, i_d=1}^{n_1, n_2, \dots, n_d} \in \mathbb{F}^{n_1 \times \dots \times n_d}$ 

Tensors

- ▶ 3-rd order **tensor**:  $\mathcal{T} = [\mathcal{T}_{ijk}]_{i,j,k=1}^{I,J,K} \in \mathbb{F}^{I \times J \times K}$
- ▶ tensor (outer) product:  $T = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$ :  $T_{ijk} = a_i b_j c_k$

Examples:

$$\mathcal{T} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} 0 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} 2 \\ -3 \end{bmatrix}$$
 
$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 0 & 2 \\ 0 & 2 \end{bmatrix},$$
 
$$\mathcal{T}_{:,:,2} = \begin{bmatrix} 0 & -3 \\ 0 & -3 \end{bmatrix},$$

 $\blacktriangleright \ \mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_2 \otimes \mathbf{e}_2 \quad \text{diagonal tensor}$ 

$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix},$$

# Some history

- late 1800s-early 1900s: algebraic geometry (Sylvester, Terracini, Segre, ...)
- ▶ 1927 (Hitchcock): introduced tensor rank
- Multiway models in psychometrics: Cattell (1940s), Tucker (1960s), Harshman (1970s)
- Popular models in chemometrics (1980s)
- Theory of complexity: Strassen (1980s)
- Signal processing: Comon (1990s)

Tensors

# Some references

Modern references (tensor decompositions):

Tensors

- ▶ [Kolda, Bader, 2009]: generic entry reference
- [Comon, 2009, 2014]: focus on CPD and its properties
- ▶ [Landsberg, 2012]: algebraic viewpoint
- [Grasedyck et al, 2013]: focus on approximation, scientific computing
- ▶ [Sidiropoulos et al, 2017]: more recent overview on uniqueness
- [Cichocki et al, 2016]: book on tensor networks

## Reminder: matrix rank

•  $\stackrel{\text{def}}{=}$  smallest r such that **X** can be factorized as

Tensors

$$m \boxed{\mathbf{X}} = m \boxed{\mathbf{A}}^{r} \boxed{\mathbf{B}}$$

 $\blacktriangleright$  = minimal r such that  $\mathbf{X}$  can be decomposed as

$$\mathbf{X} = \mathbf{a}_1^{\mathsf{T}} + \cdots + \mathbf{a}_r^{\mathsf{T}}$$

Generalization to tensors leads to two different versions of rank!

- 1. Multilinear rank (Tucker) factorization
- 2. Tensor rank (CPD) decomposition
- different decompositions for different purpose
- most other decompositions are combinations of CP and Tucker



### Matrix LRA

#### Tensors

### Factorization

(CP) decomposition

CP approximation

#### Extensions



# Generalizing matrix rank #1: multilinear rank

Matrix rank  $\stackrel{\text{def}}{=}$  dimension of column or row span:



For tensors:



tuple of (different) multilinear ranks  $(R_1, \ldots, R_d)$ ,  $R_k = \operatorname{rank}(\mathbf{Y}^{(k)})$ 

Matrix LRA Tensors

# SVDs of the unfoldings

singular values of  $\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}, \mathbf{Y}^{(3)}$ :

 $\underbrace{\mathcal{Y} \in \mathbb{R}^{610 \times 340 \times 103}}_{\text{hyperspectral image}}$ 



Pavia University hyperspectral dataset

singular values: not necessarily same distribution

Factorization

 interconnected and have well-behaved geometry [Hackbusch, Kressner, Uschmajew, 2017], [Krämer, 2019] Towards factorization: tensor/matrix product *k*-th mode contraction :

Factorization

• with  $\mathbf{M} \in \mathbb{F}^{m \times n_k}$ :  $(\mathcal{Y} \bullet_k \mathbf{M})_{i_1 \dots i_d} \stackrel{\text{def}}{=} \sum_{j=1}^{n_k} \mathcal{Y}_{i_1 \dots i_{k-1} j i_{k+1} \dots i_d} M_{i_k, j}$ For  $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$ 



For matrices (d = 2):

$$\mathbf{Y} \bullet_1 \mathbf{M} = \mathbf{M} \mathbf{Y}, \quad \mathbf{Y} \bullet_2 \mathbf{M} = \mathbf{Y} \mathbf{M}^\mathsf{T}$$

For tensors:

$$(\mathcal{Y} \bullet_k \mathbf{M})^{(k)} = \mathbf{M} \mathbf{Y}^{(k)}$$

multiplication of k-th unfolding on the left

Multilinear rank and factorization For  $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$  with ML ranks  $(R_1, R_2, R_3)$ :

Factorization

Tucker factorization with factors  $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$  (core tensor),  $\mathbf{U} \in \mathbb{R}^{I \times R_1}$ ,  $\mathbf{V} \in \mathbb{R}^{J \times R_2}$ ,  $\mathbf{W} \in \mathbb{R}^{K \times R_3}$ 



 ▶ non-unique (as in the matrix case): Ũ = UQ
 ▶ in general (random 𝒴), (R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>) = (min(I, JK), min(J, IK), min(K, IJ)) Higher-order SVD

Factorization



 $\blacktriangleright \text{ Compute } \mathcal{G}^{(SVD)} = \mathcal{Y} \bullet_1 (\mathbf{U}^{(SVD)})^{\mathsf{T}} \bullet_2 (\mathbf{V}^{(SVD)})^{\mathsf{T}} \bullet_3 (\mathbf{W}^{(SVD)})^{\mathsf{T}}$ 

$$\mathsf{HOSVD} \quad \boxed{\mathcal{Y} = \mathcal{G}^{(SVD)} \bullet_1 \mathbf{U}^{(SVD)} \bullet_2 \mathbf{V}^{(SVD)} \bullet_3 \mathbf{W}^{(SVD)}}$$

Best low multilinear approximation

Compute best  $(R_1, R_2, R_3)$ -Tucker approximation:

Factorization



a good (suboptimal) solution is given by truncating HOSVD:

$$\widehat{\mathbf{U}} = \mathbf{U}_{:,1:R_1}^{(SVD)}, \widehat{\mathbf{V}} = \mathbf{V}_{:,1:R_2}^{(SVD)}, \widehat{\mathbf{W}} = \mathbf{W}_{:,1:R_3}^{(SVD)}, \widehat{\mathcal{G}} = \mathcal{G}_{1:R_1,1:R_2,1:R_3}^{(SVD)}$$

- HOOI: alternating minimization over Û, Ŷ, Ŵ
   [De Lathauwer, De Moor, Vandewalle, 2000]
- optimization over the (R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>)-rank manifold [Kressner, Steinlechner, Vandereycken, 2013], [Kasai, Mishra, 2016]

Tucker approximation: very useful for compression, completion tasks

## Overview

### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

Extensions

Tensor rank and CPDRank-1 tensor:  $\mathcal{T} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$ :  $\mathcal{T}_{ijk} = a_i b_j c_k$ 

(Canonical) Polyadic Decomposition — sum of R rank-one tensors:

(CP) decomposition

$$\mathcal{T} = \sum_{k=1}^{R} \mathbf{a}_k \otimes \mathbf{b}_k \otimes \mathbf{c}_k, \qquad \overbrace{\mathcal{T}} = \mathbf{a}_1 \begin{vmatrix} \mathbf{c}_1 \\ \mathbf{b}_1 \end{vmatrix} + \cdots + \mathbf{c}_{R \land \mathbf{c}_R \land$$

(CP) tensor rank: rank  $(\mathcal{T}) \stackrel{\text{def}}{=} \text{minimal such } R$ 

Earlier names: CANDECOMP/PARAFAC

Example:

$$\mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_2 + \mathbf{e}_1 \otimes \mathbf{e}_2 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1$$

$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$



(CP) decomposition

$$\overbrace{\mathcal{T}}^{\mathbf{c}_{1}} = \frac{\mathbf{c}_{1}}{\mathbf{a}_{1}} + \cdots + \frac{\mathbf{c}_{R}}{\mathbf{a}_{R}}$$

Relation with multilinear ranks:



 $\max(R_1, R_2, R_3) \le R \le \min(R_2 R_3, R_1 R_3, R_1 R_2)$ 

- ▶ NP-hard (to compute exact rank): [Hillar, Lim, 2013]
- But has many nice properties and applications

$$\begin{bmatrix} \mathbf{M} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{b}_1^\mathsf{T} \\ \mathbf{a}_1 & + \cdots + & \mathbf{a}_R \\ \mathbf{a}_1 & \mathbf{a}_1 \end{bmatrix}$$

(CP) decomposition CP approximation

$$\begin{bmatrix} \mathbf{M}_{1} \\ \mathbf{M}_{1} \end{bmatrix} \stackrel{\mathbf{C}_{1,1}}{=} \mathbf{b}_{1}^{\mathsf{T}} \stackrel{\mathbf{C}_{R,1}}{=} \mathbf{b}_{R}^{\mathsf{T}} \stackrel{\mathbf{C}_$$

(CP) decomposition CP approximation

33/63

Matrix LRA Tensors Factorization (CP) decomposition CP approximation



fluorescence spectroscopy

rank = # components in a mixture emission/excitation matrix N experiments, different concentrations

Tensors Factorization (CP) decomposition CP approximation Extensions



### Tensor rank in applications

Application area	tensor rank $R$	
(blind source separation)	# of sources	
independent component analysis		
multiway factor analysis	# of components	
(spectroscopy, chemometrics,)		
antenna array processing	# of transmitters	
	:	
•	•	



Matrix LRA Tensors

## Essential uniqueness of a CPD

(CP) decomposition



up to permutations and rescaling  $(\mathbf{a}'_k = \alpha \mathbf{a}_k, \mathbf{b}'_k = \beta \mathbf{b}_k, \mathbf{c}'_k = \frac{1}{\alpha\beta} \mathbf{c}_k)$ 

#### Examples:

unique: 
$$\mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_2 \otimes \mathbf{e}_2$$
  
 $\mathcal{T}_{:,:,1} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix},$   
non-unique:  $\mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_2 + \mathbf{e}_1 \otimes \mathbf{e}_2 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_1 \otimes \mathbf{e}_2$ 

$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

Third-order tensors: Kruskal sufficient condition

Tensors Factorization (CP) decomposition CP approximation

Shorthand notation:  $\llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket := \sum_{k=1}^{R} \mathbf{a}_k \otimes \mathbf{b}_k \otimes \mathbf{c}_k$ factor matrices  $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_R]$ ,  $\mathbf{B} = [\mathbf{b}_1 \cdots \mathbf{b}_R]$ ,  $\mathbf{C} = [\mathbf{c}_1 \cdots \mathbf{c}_R]$ 

**Theorem** [Kruskal, 1978]. The decomposition  $I \xrightarrow{K} J = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$  is unique if

$$\operatorname{kr}(\mathbf{A}) + \operatorname{kr}(\mathbf{B}) + \operatorname{kr}(\mathbf{C}) \ge 2r + 2,$$

 $\underbrace{\operatorname{kr}(\mathbf{A})}_{\text{Kruskal rank}} \stackrel{\text{def}}{=} \max k \text{ such that any } k \text{ columns are linearly independent.}$ 

Example (my favourite tensor):

• K = 2 (2 frontal slices)

▶ A, B — full column rank, C — with non-collinear columns  $2R + 2 = 2R + 2 \Rightarrow$  unique decomposition

## Real vs. complex rank

(CP) decomposition

In general, they are different:

 $\mathsf{rank}_\mathbb{C}(\mathcal{T}) \leq \mathsf{rank}_\mathbb{R}(\mathcal{T})$ 

Example ( [Kruskal, 1983], but also [Sylvester, 1851]):

$$\begin{split} \mathcal{T}_{:,:,1} &= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}, \\ \mathcal{T} &= \frac{i}{2} \left( \begin{bmatrix} 1 \\ -i \end{bmatrix} \otimes \begin{bmatrix} 1 \\ -i \end{bmatrix} \otimes \begin{bmatrix} 1 \\ -i \end{bmatrix} - \begin{bmatrix} 1 \\ i \end{bmatrix} \otimes \begin{bmatrix} 1 \\ i \end{bmatrix} \otimes \begin{bmatrix} 1 \\ i \end{bmatrix} \right) \\ \Rightarrow \operatorname{rank}_{\mathbb{C}}(\mathcal{T}) &= \operatorname{rank}_{S,\mathbb{C}}(\mathcal{T}) = 2, \text{ but } \operatorname{rank}_{\mathbb{R}}(\mathcal{T}) = \operatorname{rank}_{S,\mathbb{R}}(\mathcal{T}) = 3 \end{split}$$

Symmetric tensor rank vs. tensor rank

(CP) decomposition

• 
$$S_n^d \stackrel{\text{def}}{=}$$
 vector space of symmetric tensors

Symmetric rank: rank $_S(\mathcal{T}) \stackrel{\mathrm{def}}{=} \min$  and r such that

$$\mathcal{T} = \sum_{k=1}^{r} \lambda_k \mathbf{a}_k \otimes \cdots \otimes \mathbf{a}_k, \qquad \mathbf{a}_k \in \mathbb{F}^n$$

obviously  $\mathsf{rank}(\mathcal{T}) \leq \mathsf{rank}_S(\mathcal{T})$ 

Comon's conjecture:  $\forall \mathcal{T} \in S_n^d$ ,  $\operatorname{rank}(\mathcal{T}) = \operatorname{rank}_S(\mathcal{T})$ 

true for

- matrices d = 2
- ranks smaller than order/dimension [Friedland, 2017], [Zhang, Huang, Qi, 2017]
- generically for small ranks [Lim, Qi, 2020],
- counterexample by [Shitov, 2017]: n = 800, d = 3, rank( $\mathcal{T}$ ) = 903



ation Extens

Matrix LRA Tensors Factorization (CP) decomposition CP approximation Extension

### Maximal and generic ranks



# Maximal rank

(CP) decomposition

### $r_{max} \stackrel{\text{def}}{=} \text{maximal possible rank}$

- different for symmetric/non-symmetric
- $\blacktriangleright$  different for  ${\mathbb R}$  and  ${\mathbb C}$

### (Very few) known cases

tensor space	dimension	r <sub>max</sub>	reference
$\mathbb{F}^{n \times n}$	$n^2$	n	matrices
$\mathbb{C}^{2 \times n \times n}$	$2n^2$	$\lfloor \frac{3n}{2} \rfloor$	[Grigoriev, 1978], [Ja'Ja', 1978]
$\mathbb{F}^{n  imes n  imes n}$	$n^3$	$\leq \frac{(n+1)n}{2}$	[Atkinson, Stephens, 1979]

# What if we draw tensors randomly?

(CP) decomposition

$$\overbrace{\quad }^{\mathsf{T}} \qquad \longrightarrow \mathsf{rank}\left(\mathcal{T}\right) = ?$$

random = from absolutely continuous probability distribution

- important in practice (noise, numerical errors)
- often easier to find (than  $r_{max}$ )



### Complex tensors: generic rank

With probablility 1 a complex tensor has rank  $r_{gen}$  (generic rank)

(i.e. other tensors have measure zero)



▶ matrices (ℂ<sup>n×n</sup>): r<sub>max</sub> = r<sub>gen</sub> = n
 ▶ cubic tensors (ℂ<sup>n×n×n</sup>, n > 3) [Lickteig, 1985]

$$r_{gen} = \left\lceil \frac{n^3}{3n-2} \right\rceil \approx \frac{n^2}{3}$$

## Generic ranks

(CP) decomposition

Symmetric tensors  $(S_n^d, d \ge 3)$  [Alexander, Hirschowitz, 1996]

$$r_{gen} = \left\lceil \frac{\binom{n+d-1}{n-1}}{n} \right\rceil$$

with few exceptions: (d, n) = (3, 5) or  $d = 4, n \in \{3, 4, 5\}$  $\blacktriangleright$  general case ( $\mathbb{C}^{I_1 \times \ldots \times I_d}$ ,  $d \ge 3$ ):

$$r_{gen} \stackrel{?}{=} \left[ \frac{I_1 \cdots I_d}{I_1 + \cdots + I_d - d + 1} \right], \text{ with few exceptions}$$

▶ conjectured [Abo, Ottaviani, Peterson, 2009]
 ▶ computer proof [Chiantini, Ottaviani, Vannieuwenhoven, 2014]
 ▶ [Blekherman, Teitler, 2014]: r<sub>max</sub> ≤ 2r<sub>gen</sub>

Proofs: algebraic geometry (dimensions of secant varieties)

# Real tensors: typical ranks

(CP) decomposition

For real tensors, several typical ranks may appear with nonzero probability.



Example [Bergqvist, 2013]:  $T \in \mathbb{R}^{2 \times 2 \times 2}$  with i.i.d. Gaussian elements has:

rank 2 with probability π/4
 rank 3 with probability 1 - π/4

44 / 63

## Some numerical consequences

(CP) decomposition

1. noise, numerical errors  $\Rightarrow$  rank $(\mathcal{T}) = r_{gen}$  (or a typical rank in  $\mathbb{R}$ )

2. very difficult to find tensors with higher ranks:

If we generate

$$\underbrace{\mathcal{T}}_{\mathbf{a}_1 \mid \mathbf{b}_1} = \underbrace{\mathbf{c}_1 / \mathbf{b}_1}_{\mathbf{a}_1 \mid \mathbf{b}_1} + \cdots + \underbrace{\mathbf{c}_r / \mathbf{b}_r}_{\mathbf{a}_r \mid \mathbf{b}_r}$$

with  $\mathbf{a}_k$ ,  $\mathbf{b}_k$ ,  $\mathbf{c}_k$  random, has

$$\mathsf{rank}(\mathcal{T}) = \begin{cases} r, & r \leq r_{gen}, \\ r_{gen}, & r_{gen} < r \leq r_{max} \end{cases}$$

**Example**.  $n \times n \times n$  tensors with rank  $\left\lceil \frac{n^3}{3n-2} \right\rceil < r \leq \frac{(n+1)n}{2}$ .
# Generic uniqueness (identifiability)

For a fixed rank r: whether "almost all" decompositions are unique

- Kruskal-type conditions give weak bounds
- Study the properties secant algebraic variety (σ<sub>r</sub> <sup>def</sup>=: (Zariski) closure of tensors of rank r)
- $\blacktriangleright$  Generic uniqueness: uniqueness for all tensors in  $\sigma_r$  except a set of Lebesgue measure 0

Most recent results:

- 1. [Chiantini, Ottaviani, 2012]: CPD of  $\mathcal{T} \in \mathbb{C}^{I \times J \times K}$ , with  $I \ge J \ge K$  is generically unique if  $r \le 2^{\lfloor \log_2 J \rfloor + \lfloor \log_2 K \rfloor 2}$ .
- 2. [Chiantini, Ottaviani, Vannieuwenhoven, 2014] (computer proof): complex identifiability holds for all subgeneric ranks  $r < r_{gen} = \left\lceil \frac{IJK}{I+J+K-2} \right\rceil$ no identifiability for  $r > r_{gen}$
- 3. [Qi, Comon, Lim, 2016], [Chiantini, Ottaviani, Vanniueuwenhoven, 2017]: all identifiability results are valid for real-valued tensors

## Overview

### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

### CP vs. Tucker

CP approximation

• CPD: 
$$(I + J + K - 2)R$$

$$\overbrace{\mathcal{T}}^{\mathbf{c}_{1}} \approx_{\mathbf{a}_{1}}^{\mathbf{c}_{1}} + \cdots +_{\mathbf{a}_{R}}^{\mathbf{c}_{R}} \mathbf{b}_{R}$$

▶ Tucker:  $IR_1 + JR_2 + KR_3 + R_1R_2R_3 - R_1^2 - R_2^2 - R_3^2$ 





CP approximation

$$\overbrace{\mathcal{T}}^{\mathbf{c}_{1}} \approx_{\mathbf{a}_{1}}^{\mathbf{c}_{1}} + \cdots +_{\mathbf{a}_{R}}^{\mathbf{c}_{R}} \mathbf{b}_{R}$$

Best r-rank approximation (r > 1) may not exist

$$\mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_2 + \mathbf{e}_1 \otimes \mathbf{e}_2 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1$$
$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

rank  $\mathcal{T} = 3$ , but approximated by rank-2 tensor to any accuracy

$$\mathcal{T} = \frac{1}{2\varepsilon} (\mathbf{e}_1 + \varepsilon \mathbf{e}_2)^{\otimes 3} - \frac{1}{2\varepsilon} (\mathbf{e}_1 - \varepsilon \mathbf{e}_2)^{\otimes 3} + \mathcal{O}(\varepsilon)$$

Set of rank- $\leq r$  tensors is not closed for r > 1

Rank-one tensor approximation

CP approximation

 $\min_{\mathbf{a},\mathbf{b},\mathbf{c}} \| \mathcal{T} - \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} \|_F^2$ 

- well-posed (minimum exists)
- block-coordinate descent (ALS, non-symmetric power method) converges globally (to a stationary point) [Uschmajew, 2015]
- related to the notion of singular vectors/eigenvectors of a tensor [Qi, Luo, 2017]
- number of stationary points is known [Freidland, Ottaviani, 2014]

# Successive approximation (deflation)

CP approximation

Subtracting rank-one approximation may increase tensor rank:

$$\mathcal{X}_{:,:,1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{X}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

 $\operatorname{rank}(\mathcal{X})=2$ , but  $\mathcal{X}-\widehat{\mathcal{X}}_1=\mathcal{T}$ ,  $\operatorname{rank}(\mathcal{T})=3$ 

 $\mathcal{T} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_2 + \mathbf{e}_1 \otimes \mathbf{e}_2 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1$ 

$$\mathcal{T}_{:,:,1} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathcal{T}_{:,:,2} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

But when does it work then?

# Successive rank-1 approximation

CP approximation

Orthogonally decomposable tensors [Zhang, Golub, 2001]:

$$\mathcal{T} = \sum_{k=1}^{R} \mathbf{a}_k \otimes \mathbf{b}_k \otimes \mathbf{c}_k, \quad \mathbf{a}_k ot \mathbf{a}_j, \quad \mathbf{b}_k ot \mathbf{b}_j, \quad \mathbf{c}_k ot \mathbf{c}_j,$$

The successive best rank-1 approximation returns the components in the sum

Cyclic (sequential) rank-one approximation [da Silva, Comon, de Almeida, 2015]:

$$\min \|\mathcal{T} - (\widehat{\mathcal{T}}_1 + \ldots + \widehat{\mathcal{T}}_R)\|$$
 subject to  $\operatorname{rank}(\widehat{\mathcal{T}}_k) = 1$ .

•  $\mathcal{T} \in S_n^d$  — symmetric tensors,



**CP** approximation

- Best non-symmetric approximation can be chosen symmetric (and is unique a.s.) [Friedland, Ottaviani, 2014]
- $\blacktriangleright$  Best symmetric approximation  $\leftrightarrow$  maximization of a polynomial

$$\min_{\lambda,\mathbf{v}} \|\mathcal{T} - \lambda \mathbf{v} \otimes \cdots \otimes \mathbf{v}\|_F^2 \leftrightarrow \max_{\|\mathbf{v}\|_2 = 1} |\mathcal{T} \bullet_1 \mathbf{v} \cdots \bullet_d \mathbf{v}|$$

stationary points: eigenvectors of the tensor

$$\mathcal{T} \cdot \mathbf{v}^{d-1} = \mu \mathbf{v}$$

In total, (d-1)<sup>n</sup>-1/(d-2) (complex) eigenvectors [Cartwright, Sturmfels, 2013]
 There are cases when the power method diverges [Chen, Saad, 2009]
 Orthogonally decomposable tensors: power method converges, deflation works decomposition [Anandkumar et al, 2013], [Robeva.2016]

# Algorithms for CP approximation

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \|\mathcal{Y} - [\![\mathbf{A},\mathbf{B},\mathbf{C}]\!]\|_F^2 + \dots$$

Constraints/regularization:

- Orthogonality: [Comon, 1993], [Robeva, 2016]
- Coherence: [Lim, Comon, 2014]
- Nonnegativity: [Qi, Comon, Lim, 2016]

Algorithms for CP approximation

- Alternating minimization (ALS, AO-ADMM) (Global) convergence properties in the regularized case [Xu, Yin, 2013]
- Nonlinear least squares
- Riemannian optimization
- Algebraic algorithms
  - Generalized eigenvalue decomposition (non-symmetric tensors)
  - Structured matrix approximation (symmetric tensors)

approximation

Tensor diagonalization as orthogonal approximation By duality:

$$\max_{\mathbf{Q}\in\mathcal{O}_{n}} \|\operatorname{diag}\mathcal{A}\bullet_{1}\mathbf{Q}^{\mathsf{T}}\cdots\bullet_{d}\mathbf{Q}^{\mathsf{T}}\|_{2}^{2}$$

$$= \|\mathcal{A}\|_{F}^{2} - \min_{\substack{\mathbf{Q}\in\mathcal{O}_{n}\\ \mathbf{u}_{1}\cdots\cdot\mathbf{u}_{n}}} \|\mathcal{A}-\sum_{k=1}^{n}\mu_{k}\mathbf{u}_{k}\otimes\cdots\otimes\mathbf{u}_{k}\|_{F}^{2}$$
best *n*-rank symmetric orthogonal approximation

Euclidean distance to odeco [Robeva, 2016] variety



▶ Best non-symmetric approximation ≠ best symmetric approximation [Li, U., Comon, 2019] (a variant of Comon's conjecture is false)

approximation

**CP** approximation

 $\max_{\mathbf{Q}\in\mathcal{O}_n}f(\mathbf{Q})\quad\text{or}\quad\max_{\mathbf{U}\in\mathcal{U}_n}f(\mathbf{U}),\quad f\text{ is a low-order polynomial}$ 

Algebraic orthogonal tensor decompositions:

- Deflation (successive rank-one approximation) [Delfosse, Loubaton, 1995], [Anandkumar, 2013], [Robeva, 2016]
- ▶ EVD of tensor slice(s) [De Lathauwer, 2006], [Kolda, 2015]

Optimization on the manifold:

- ▶ Riemannian optimization (CG, SD, BFGS, RTR) [Absil et al., 2008]
- Jacobi-type algorithms [Comon, 1993], [Li, U., Comon, 2020] (convergence)

What do we know about CP approximation?

**CP** approximation

Given  $\mathcal{T} = [\![\mathbf{A},\mathbf{B},\mathbf{C}]\!] + \mathcal{E}$  and

$$(\widehat{\mathbf{A}}^*, \widehat{\mathbf{B}}^*, \widehat{\mathbf{C}}^*) = \arg\min_{\widehat{\mathbf{A}}, \widehat{\mathbf{B}}, \widehat{\mathbf{C}}} \|\mathcal{T} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]\|,$$

- Approaches using random matrix theory [Goulart, Couillet, Comon, 2022] (mostly rank-1)
- Perturbation bounds for some algebraic algorithms [Evert, De Lathauwer, 2022] (very small ranks)
- Uniqueness of best approximation in a small neighbrhood [Friedland, Stawiska, 2016] (non-constructive)

## Overview

### Matrix LRA

#### Tensors

Factorization

(CP) decomposition

CP approximation

# Coupled factorizations

Factors of different tensor/matrix decompositions may be shared:



[Acar, Kolda, Dunlavy, 2011]

## Joint CPD of symmetric tensors

ICA model:  $\mathbf{x} = A\mathbf{s}, \quad A = \begin{bmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_r \end{bmatrix} \in \mathbb{R}^{n \times r}$ ,

cumulants of  $\mathbf{x}$  up to order d:

$$\begin{aligned}
\mathcal{C}_{\mathbf{x}}^{(1)} &= c_{1,1}\mathbf{a}_{1} + \dots + c_{1,r}\mathbf{a}_{r}, \\
\mathcal{C}_{\mathbf{x}}^{(2)} &= c_{2,1}\mathbf{a}_{1} \otimes \mathbf{a}_{1} + \dots + c_{2,r}\mathbf{a}_{r} \otimes \mathbf{a}_{r}, \\
&\vdots \\
\mathcal{C}_{\mathbf{x}}^{(d)} &= c_{d,1}\mathbf{a}_{1} \otimes \dots \otimes \mathbf{a}_{1} + \dots + c_{d,r}\mathbf{a}_{r} \otimes \dots \otimes \mathbf{a}_{r},
\end{aligned}$$
(1)

where  $c_{j,k}$  is the *j*-th cumulant of  $s_k$ .

**Problem**. Given  $C_{\mathbf{x}}^{(j)}$ , find  $\mathbf{a}_k$  and  $c_{j,k}$  (diagonalize all the cumulants simultaneously)

Matrix LRA Tensors Factorization (CP) decomposition CP approximation E

### Sum of Tucker: block-term decomposition

[De Lathauwer, 2008]: For fixed  $(R_1, R_2, R_3)$ :

$$\mathcal{T} = \sum_{k=1}^{r} \mathcal{G}_k \bullet_1 \mathbf{U}_k \bullet_2 \mathbf{V}_k \bullet_3 \mathbf{W}_k$$

each term in the sum has ML-rank  $(R_1, R_2, R_3)$ :



Special case:  $(R_1, R_2, R_3) = (L, L, 1)$ : very useful in signal processing, see e.g. [Goulart, et al. 2020]

## Additive decompositions

X-rank (join) decomposition [Zak, 2004], [Landsberg, 2012], [Comon, Qi, U., 2017] ( sparse algebraic decomposition)

 $\mathbf{p} = \mathbf{x}_1 + \dots + \mathbf{x}_r, \quad \mathbf{x}_k \in \widehat{X}$ 

- $A \stackrel{\text{def}}{=}$  ambient tensor space (e.g.  $A = \mathbb{C}^{I \times J \times K}$ )
- Â def = (variety of "simple" ' terms)

   (e.g., rank-one X̂ = {a ⊗ b ⊗ c})
- ► study the properties of secant varieties  $\sigma_r(\hat{X})$



# Approximation: higher-order tensors

• Tucker:  $O(dIR + R^d)$ : curse of dimensionality



tensor trains, hierarchical Tucker: linear-in-d storage complexity



Decompositions: flexible/multilayer [Harshman, Lundy, 1996], [Roald et al, 2022] CPD (PARAFAC)

$$\mathcal{T}_{:,:,k} = \prod_{I \in \mathbf{A}} \frac{R}{\mathbf{A}} \cdot \underbrace{\mathbf{D}_{\mathbf{C}}^{(k)}}_{\mathbf{C}} \cdot R \underbrace{\mathbf{B}^{\mathsf{T}}}_{\mathbf{B}^{\mathsf{T}}}, \quad k = 1, \dots, K$$

PARAFAC-2

$$\mathcal{T}_{:,:,k} = \prod_{I \in \mathbf{A}} \frac{R}{\mathbf{A}} \cdot \underbrace{\mathbf{D}_{\mathbf{C}}^{(k)}}_{\mathbf{C}} \cdot R \underbrace{\mathbf{B}_{k}^{\mathsf{T}}}_{\mathbf{B}_{k}^{\mathsf{T}}}, \quad k = 1, \dots, K$$

ParaTuck-2

$$\mathcal{T}_{:,:,k} = \prod_{I \in \mathbf{A}} \mathbf{A} \cdot \mathbf{D}_{\mathbf{C}}^{(k)} \cdot R \mathbf{F} \cdot \mathbf{D}_{\mathbf{H}}^{(k)} \cdot S \mathbf{B}^{\mathsf{T}}, \quad k = 1, \dots, K$$

Share the uniqueness features of CPD!

# Conclusion

Which tensor decomposition (format) to use?

- ▶ factorization (Tucker, tensor trains) compression
- decomposition (CPD, BTD) identification of components

Challenges (personal choice):

- Guarantees (approximation bounds) on CP approximation
- Multi-layer/multilevel tensor decompositions (e.g., ParaTuck-2)

# Conclusion

Which tensor decomposition (format) to use?

- ▶ factorization (Tucker, tensor trains) compression
- decomposition (CPD, BTD) identification of components

Challenges (personal choice):

- Guarantees (approximation bounds) on CP approximation
- Multi-layer/multilevel tensor decompositions (e.g., ParaTuck-2)

Advertisement (not mine ):

- A number of postdoc/PhD positions in the SiMul team (Nancy) https://cran-simul.github.io
- Summer school on low-rank approximation (23-29.06.24) in Peyresq Organized by N. Gillis and J. Cohen

Thank you!