

Tensor models & applications

J. Henrique de M. Goulart

Toulouse INP / IRIT

`henrique.goulart@irit.fr`

Peyresq – June 2024



About this mini-course

Objectives & method:

1. Provide an overview of the **main concepts & results** pertaining to tensors & their decompositions in **signal processing & machine learning**.
2. Emphasis on **intuitions** rather than formal proofs (but still keep some).
3. Rely on application examples from these domains to **motivate & illustrate** the introduced tools.

Disclaimer:

- This is a very broad subject with a **vast literature**, spanning several communities (SP, ML, chemometrics, numerical analysis, physics, ...).
- This mini-course is **my personal take** on it—hence partial and biased.

Prerequisites & handout

Prerequisites:

- linear algebra
- basic probability & statistics
- basic optimization

Supporting material: An electronic [handout](https://www.irit.fr/~Henrique.Goulart/talks/) is provided on my webpage:

`https://www.irit.fr/~Henrique.Goulart/talks/`

It contains:

- a summary of important identities and properties
- some proposed exercices
- many bibliographical pointers

Agenda

- 0. About this mini-course
- 1. Why care about tensors?
- 2. A jungle of tensor models
- 3. A mosaic of uniqueness results
- 4. Rank-1 approximation, tensor spectrum and power iteration
- 5. Low-rank approximation, in several flavors
- 6. Tensor PCA & asymptotic MLE performance

Why care about tensors?



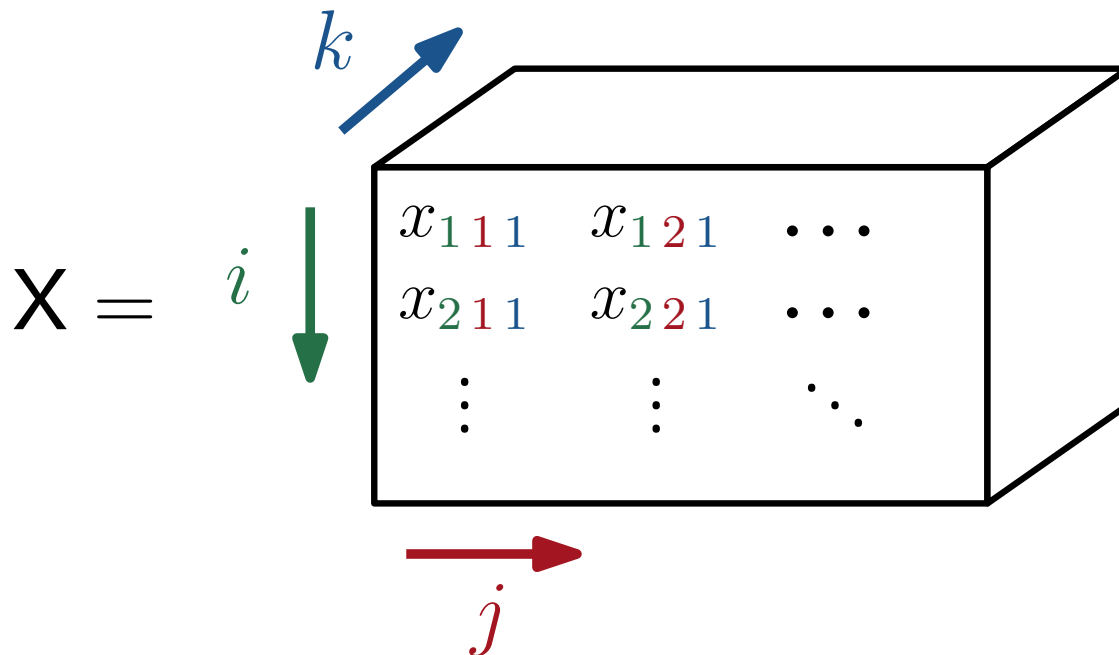
A first definition

Tensors are (for now) simply **multi-way arrays** (a.k.a. hypermatrices)

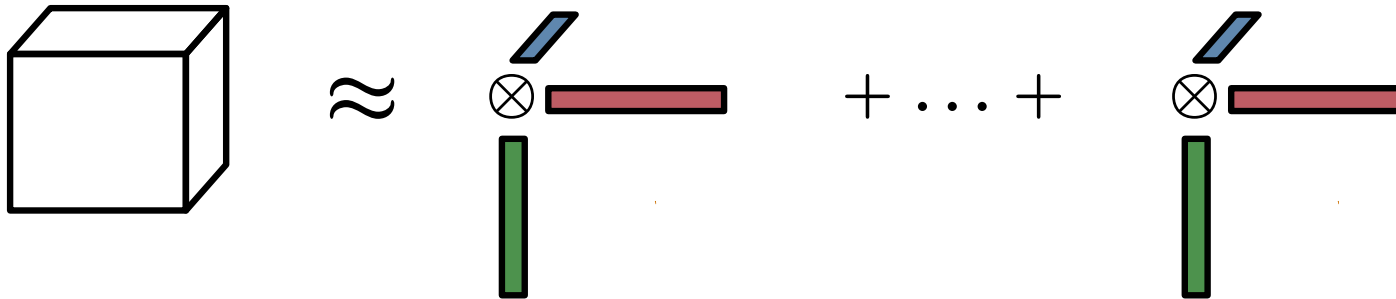
$$\mathbf{X} \in \mathbb{R}^{N_1 \times \cdots \times N_d}, \quad (\mathbf{X})_{i_1 \dots i_d} = x_{i_1 \dots i_d},$$

where d is the **order** of \mathbf{X} .

For $d = 3$, we write $\mathbf{X} = (\mathbf{X})_{ijk}$. We often say that \mathbf{X} has **d modes**.

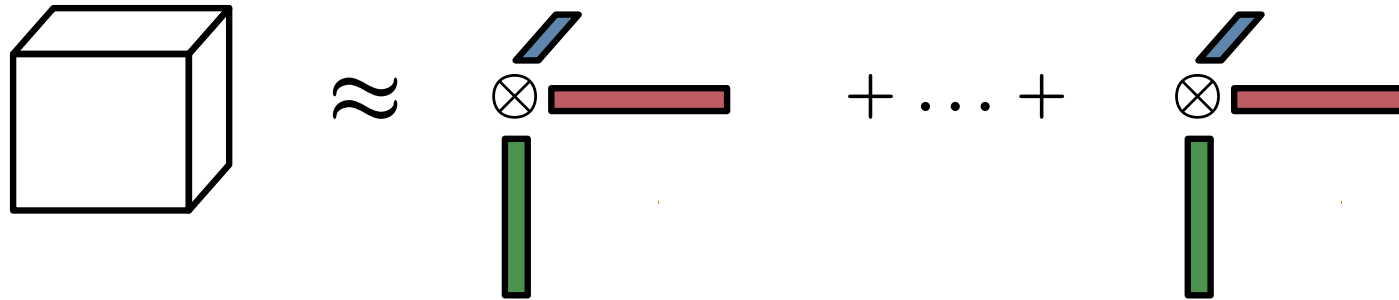


Numerous applications involve data (or functions thereof) or other objects that can be modelled by tensors having a **low-dimensional (low-rank) structure**.



Tensor models

Numerous applications involve data (or functions thereof) or other objects that can be modelled by tensors having a **low-dimensional (low-rank) structure**.



Such **tensor models** are useful for many purposes falling into two categories:

1. Information extraction
2. Complexity reduction

Let's see a few examples.

Example #1: estimating Gaussian mixtures

Take a p th-dim spherical Gaussian mixture model (GMM):

$$\mathbf{X} \sim \sum_{r=1}^R \pi_r \mathcal{N}(\boldsymbol{\mu}_r, \sigma_r^2 \mathbf{I}).$$

Goal: Given observations $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^N$, estimate $\{(\pi_r, \boldsymbol{\mu}_r, \sigma_r^2)\}_{r=1}^R$.

Classical approach: expectation-minimization (EM) algorithm¹—sensitive w.r.t. initialization, may converge very slowly.²

1: Dempster & al., 1977, 2: Park & Ozeki, 2009

Example #1: estimating Gaussian mixtures

Take a p -th-dim spherical Gaussian mixture model (GMM):

$$\mathbf{X} \sim \sum_{r=1}^R \pi_r \mathcal{N}(\boldsymbol{\mu}_r, \sigma_r^2 \mathbf{I}).$$

Goal: Given observations $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^N$, estimate $\{(\pi_r, \boldsymbol{\mu}_r, \sigma_r^2)\}_{r=1}^R$.

Classical approach: expectation-minimization (EM) algorithm¹—sensitive w.r.t. initialization, may converge very slowly.²

Let's try instead the **method of moments**:

$$\begin{aligned} \mathbb{E}\{\mathbf{X}\} &= \sum_{r=1}^R \pi_r \boldsymbol{\mu}_r, \\ \mathbb{E}\{\mathbf{X}\mathbf{X}^\top\} &= \sum_{r=1}^R \pi_r \left(\boldsymbol{\mu}_r \boldsymbol{\mu}_r^\top + \sigma_r^2 \mathbf{I} \right) = \sum_{r=1}^R \pi_r \boldsymbol{\mu}_r \boldsymbol{\mu}_r^\top + \underbrace{\sum_r \pi_r \sigma_r^2}_{\bar{\sigma}^2} \mathbf{I}. \end{aligned}$$

Q: Given $\mathbb{E}\{\mathbf{X}\}$, $\mathbb{E}\{\mathbf{X}\mathbf{X}^\top\}$, can we identify the GMM parameters?

1: Dempster & al., 1977, 2: Park & Ozeki, 2009

First try: second-order moments

Letting $\tilde{\boldsymbol{\mu}}_r := \sqrt{\pi_r} \boldsymbol{\mu}_r$ and $\mathbf{M} := (\tilde{\boldsymbol{\mu}}_1 \ \dots \ \tilde{\boldsymbol{\mu}}_R)$, we can write:

$$\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} = \sum_{r=1}^R \tilde{\boldsymbol{\mu}}_r \tilde{\boldsymbol{\mu}}_r^\top + \bar{\sigma}^2 \mathbf{I} = \mathbf{M} \mathbf{M}^\top + \bar{\sigma}^2 \mathbf{I}$$

First try: second-order moments

Letting $\tilde{\boldsymbol{\mu}}_r := \sqrt{\pi_r} \boldsymbol{\mu}_r$ and $\mathbf{M} := (\tilde{\boldsymbol{\mu}}_1 \ \dots \ \tilde{\boldsymbol{\mu}}_R)$, we can write:

$$\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} = \sum_{r=1}^R \tilde{\boldsymbol{\mu}}_r \tilde{\boldsymbol{\mu}}_r^\top + \bar{\sigma}^2 \mathbf{I} = \mathbf{M} \mathbf{M}^\top + \bar{\sigma}^2 \mathbf{I}$$

Let's focus on estimating the means $\boldsymbol{\mu}_r$.

If $R < N$, then $\bar{\sigma}^2 = \lambda_{\min} \left(\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} \right)$, estimated by $\lambda_{\min} \left(\frac{1}{M} \sum_{m=1}^M \mathbf{x}_m \mathbf{x}_m^\top \right)$.

Hence, we can approximately compute

$$\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} - \bar{\sigma}^2 \mathbf{I} = \mathbf{M} \mathbf{M}^\top.$$

Can we then recover $\{\tilde{\boldsymbol{\mu}}_r\}_{r=1}^R$ by matrix decomposition?

First try: second-order moments

Letting $\tilde{\boldsymbol{\mu}}_r := \sqrt{\pi_r} \boldsymbol{\mu}_r$ and $\mathbf{M} := (\tilde{\boldsymbol{\mu}}_1 \ \dots \ \tilde{\boldsymbol{\mu}}_R)$, we can write:

$$\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} = \sum_{r=1}^R \tilde{\boldsymbol{\mu}}_r \tilde{\boldsymbol{\mu}}_r^\top + \bar{\sigma}^2 \mathbf{I} = \mathbf{M} \mathbf{M}^\top + \bar{\sigma}^2 \mathbf{I}$$

Let's focus on estimating the means $\boldsymbol{\mu}_r$.

If $R < N$, then $\bar{\sigma}^2 = \lambda_{\min} \left(\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} \right)$, estimated by $\lambda_{\min} \left(\frac{1}{M} \sum_{m=1}^M \mathbf{x}_m \mathbf{x}_m^\top \right)$.

Hence, we can approximately compute

$$\mathbb{E} \left\{ \mathbf{X} \mathbf{X}^\top \right\} - \bar{\sigma}^2 \mathbf{I} = \mathbf{M} \mathbf{M}^\top.$$

Can we then recover $\{\tilde{\boldsymbol{\mu}}_r\}_{r=1}^R$ by matrix decomposition? **No:**

$$\mathbf{M} \mathbf{M}^\top = (\mathbf{M} \mathbf{S})(\mathbf{M} \mathbf{S}^{-\top})^\top, \quad \forall \mathbf{S} \in \text{GL}_R(\mathbb{R}).$$

“too much freedom = not identifiable”

Going to higher orders

The third-order moments give $\binom{N+2}{3}$ distinct equations of the form:

$$\mathbb{E} \{X_i X_j X_k\} = \sum_{r=1}^R \pi_r \mathbb{E} \{[(\boldsymbol{\mu}_r)_i + (\mathbf{U}_r)_i] [(\boldsymbol{\mu}_r)_j + (\mathbf{U}_r)_j] [(\boldsymbol{\mu}_r)_k + (\mathbf{U}_r)_k]\}$$

where $\mathbf{U}_r \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{I})$.

Going to higher orders

The third-order moments give $\binom{N+2}{3}$ distinct equations of the form:

$$\mathbb{E} \{ X_i X_j X_k \} = \sum_{r=1}^R \pi_r \mathbb{E} \{ [(\boldsymbol{\mu}_r)_i + (\mathbf{U}_r)_i] [(\boldsymbol{\mu}_r)_j + (\mathbf{U}_r)_j] [(\boldsymbol{\mu}_r)_k + (\mathbf{U}_r)_k] \}$$

where $\mathbf{U}_r \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{I})$. Direct computation gives

$$\mathbb{E} \{ X_i X_j X_k \} = \sum_{r=1}^R \pi_r \left[(\boldsymbol{\mu}_r)_i (\boldsymbol{\mu}_r)_j (\boldsymbol{\mu}_r)_k + (\boldsymbol{\mu}_r)_i \sigma_r^2 \delta_{jk} + \right. \\ \left. (\boldsymbol{\mu}_r)_j \sigma_r^2 \delta_{ik} + (\boldsymbol{\mu}_r)_k \sigma_r^2 \delta_{ij} \right]$$

Going to higher orders

The third-order moments give $\binom{N+2}{3}$ distinct equations of the form:

$$\mathbb{E} \{X_i X_j X_k\} = \sum_{r=1}^R \pi_r \mathbb{E} \{[(\boldsymbol{\mu}_r)_i + (\mathbf{U}_r)_i] [(\boldsymbol{\mu}_r)_j + (\mathbf{U}_r)_j] [(\boldsymbol{\mu}_r)_k + (\mathbf{U}_r)_k]\}$$

where $\mathbf{U}_r \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{I})$. Direct computation gives

$$\mathbb{E} \{X_i X_j X_k\} = \sum_{r=1}^R \pi_r \left[(\boldsymbol{\mu}_r)_i (\boldsymbol{\mu}_r)_j (\boldsymbol{\mu}_r)_k + (\boldsymbol{\mu}_r)_i \sigma_r^2 \delta_{jk} + \right. \\ \left. (\boldsymbol{\mu}_r)_j \sigma_r^2 \delta_{ik} + (\boldsymbol{\mu}_r)_k \sigma_r^2 \delta_{ij} \right]$$

By introducing an appropriate symmetrization operator Sym , we get

$$\mathbb{E} \{X_i X_j X_k\} = \sum_{r=1}^R \pi_r (\boldsymbol{\mu}_r)_i (\boldsymbol{\mu}_r)_j (\boldsymbol{\mu}_r)_k + 3 \text{Sym} \left(\delta_{jk} \underbrace{\sum_{r=1}^R \pi_r \sigma_r^2 (\boldsymbol{\mu}_r)_i}_{:= (\mathbf{c})_i} \right)$$

A first tensor decomposition

Let's introduce some convenient (& natural) notation:

Def: elementary or rank-1 tensor

$$(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \in \mathbb{R}^{N_1} \times \mathbb{R}^{N_2} \times \mathbb{R}^{N_3} \mapsto \mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

$$(\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3)_{ijk} = (\mathbf{x}_1)_i (\mathbf{x}_2)_j (\mathbf{x}_3)_k$$

A first tensor decomposition

Let's introduce some convenient (& natural) notation:

Def: elementary or rank-1 tensor

$$(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \in \mathbb{R}^{N_1} \times \mathbb{R}^{N_2} \times \mathbb{R}^{N_3} \mapsto \mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

$$(\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3)_{ijk} = (\mathbf{x}_1)_i (\mathbf{x}_2)_j (\mathbf{x}_3)_k$$

This allows expressing all $\binom{N+2}{3}$ equations as:

$$\mathbb{E} \{ \mathbf{X} \otimes \mathbf{X} \otimes \mathbf{X} \} = \sum_{r=1}^R \pi_r \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r + \mathbf{S}(\mathbf{c})$$

Moreover, \mathbf{c} can be estimated from data.

A first tensor decomposition

Let's introduce some convenient (& natural) notation:

Def: elementary or rank-1 tensor

$$(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \in \mathbb{R}^{N_1} \times \mathbb{R}^{N_2} \times \mathbb{R}^{N_3} \mapsto \mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

$$(\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3)_{ijk} = (\mathbf{x}_1)_i (\mathbf{x}_2)_j (\mathbf{x}_3)_k$$

This allows expressing all $\binom{N+2}{3}$ equations as:

$$\mathbb{E} \{ \mathbf{X} \otimes \mathbf{X} \otimes \mathbf{X} \} = \sum_{r=1}^R \pi_r \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r + \mathbf{S}(\mathbf{c})$$

Moreover, \mathbf{c} can be estimated from data.

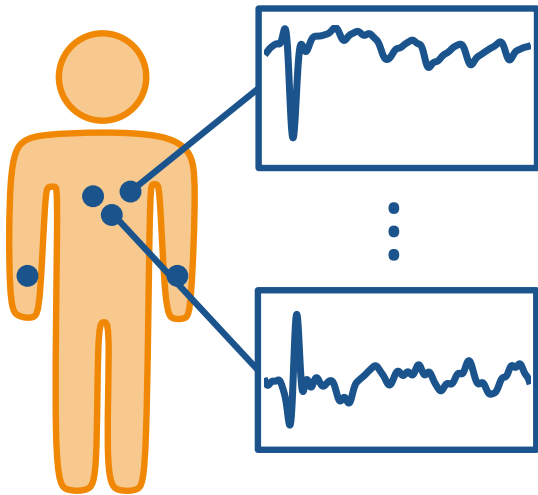
Hence, our problem \approx decomposition of the third-order moment tensor:^{1,2}

$$\mathbb{E} \{ \mathbf{X} \otimes \mathbf{X} \otimes \mathbf{X} \} - \mathbf{S}(\mathbf{c}) = \sum_{r=1}^R \pi_r \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r \otimes \boldsymbol{\mu}_r.$$

Good news: (essentially) unique decomposition under mild constraints!

1: Hsu & Kakade, 2013, 2: Anandkumar & al., 2014

Example #2: ECG signal separation



$$x_1(n) = \sum_{r=1}^R (w_r)_1 s_r(n)$$

$$\vdots$$

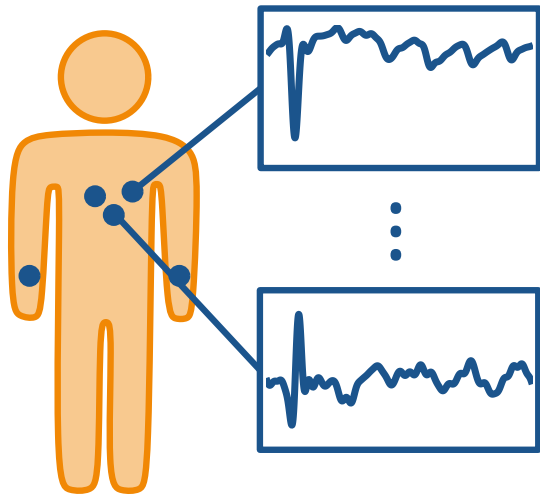
$$x_K(n) = \sum_{r=1}^R (w_r)_K s_r(n)$$

K sensors (leads)

N observations (samples)

R sources

Example #2: ECG signal separation



$$x_1(n) = \sum_{r=1}^R (w_r)_1 s_r(n)$$

K sensors (leads)

N observations (samples)

$$x_K(n) = \sum_{r=1}^R (w_r)_K s_r(n)$$

R sources

The data can be expressed as a matrix decomposition (linear mixture model):

$$\underbrace{\begin{pmatrix} x_1 & \dots & x_K \end{pmatrix}}_{\mathbf{X} \in \mathbb{R}^{N \times K}} = \underbrace{\begin{pmatrix} s_1 & \dots & s_R \end{pmatrix}}_{\mathbf{S} \in \mathbb{R}^{N \times R}} \underbrace{\begin{pmatrix} w_1 & \dots & w_R \end{pmatrix}^T}_{\mathbf{W}^T \in \mathbb{R}^{R \times K}}$$

$$= \underbrace{\mathbf{S} \mathbf{P}}_{\tilde{\mathbf{S}}} \underbrace{\mathbf{P}^{-1} \mathbf{W}^T}_{\tilde{\mathbf{W}}^T}, \quad \forall \mathbf{P} \in \text{GL}_R(\mathbb{R}).$$

This model is clearly **non-identifiable** without further constraints.

“Tensorization” by additional diversity

Idea: “Hankelize” the observed signals to add one temporal “diversity”.¹

$$\begin{aligned}
 \mathbf{x}_k \mapsto \mathbf{X}_k = & \begin{array}{c} x_k(0) \rightarrow \begin{array}{|c|c|c|c|c|c|c|c|} \hline \text{orange} & \text{white} & \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} \\ \hline \text{white} & \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} \\ \hline \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} & \text{green} \\ \hline \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} & \text{green} & \text{white} \\ \hline \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} & \text{green} & \text{white} & \text{red} \\ \hline \text{red} & \text{white} & \text{orange} & \text{green} & \text{white} & \text{red} & \text{blue} & \text{orange} \\ \hline \text{white} & \text{orange} & \text{green} & \text{white} & \text{red} & \text{blue} & \text{orange} & \text{white} \\ \hline \end{array} \\ \downarrow \\ x_k(N-1) \rightarrow \end{array} \\
 = & \sum_{r=1}^R (\mathbf{w}_r)_k \mathbf{H}_r \\
 & \quad \quad \quad \text{red} \quad \quad \quad \text{red} \\
 & \quad \quad \quad \vdots \\
 & \quad \quad \quad \mathbf{s}_r \mapsto \mathbf{H}_r
 \end{aligned}$$

“Tensorization” by additional diversity

Idea: “Hankelize” the observed signals to add one temporal “diversity”.¹

$$x_k \mapsto \mathbf{X}_k = \begin{array}{c} x_k(0) \text{---} \begin{array}{|c|c|c|c|c|c|c|} \hline \text{orange} & \text{white} & \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} \\ \hline \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} \\ \hline \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} & \text{white} & \text{orange} \\ \hline \text{orange} & \text{white} & \text{green} & \text{blue} & \text{orange} & \text{blue} & \text{red} \\ \hline \text{red} & \text{white} & \text{orange} & \text{blue} & \text{orange} & \text{blue} & \text{red} \\ \hline \text{red} & \text{white} & \text{orange} & \text{blue} & \text{orange} & \text{blue} & \text{red} \\ \hline \end{array} \text{---} x_k(N-1) \end{array} = \sum_{r=1}^R (\mathbf{w}_r)_k \mathbf{H}_r$$

$s_r \mapsto \mathbf{H}_r$

Stacking the \mathbf{X}_k as slices of $\mathbf{X} \in \mathbb{R}^{M \times M \times K}$, with $M = \frac{N+1}{2}$, we get:

$$\mathbf{X} = \begin{array}{c} \text{---} \mathbf{X}_K \\ \mathbf{X}_1 \end{array} = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{w}_r = \sum_{r=1}^R \begin{array}{c} \text{---} \mathbf{w}_r \\ \mathbf{H}_r \end{array}$$

Def: matrix-vector tensor product $(\mathbf{H} \otimes \mathbf{w})_{ijk} = (\mathbf{H})_{ik} (\mathbf{w}_r)_k$

A decomposition in low-rank blocks

Without further assumptions, the model is still non-identifiable.

But we can add a reasonable one:

“Parsimony” assumption: sources are given by

$$s_r(n) = \sum_{\ell=1}^{L_r} \alpha_{\ell,r} z_{\ell,r}^n \quad \Rightarrow \quad \text{rank } \mathbf{H}_r \leq L_r < M,$$

where $\alpha_{\ell,r}, z_{\ell,r} \in \mathbb{C}$.

A decomposition in low-rank blocks

Without further assumptions, the model is still non-identifiable.

But we can add a reasonable one:

“Parsimony” assumption: sources are given by

$$s_r(n) = \sum_{\ell=1}^{L_r} \alpha_{\ell,r} z_{\ell,r}^n \quad \Rightarrow \quad \text{rank } \mathbf{H}_r \leq L_r < M,$$

where $\alpha_{\ell,r}, z_{\ell,r} \in \mathbb{C}$.

Hence, our source separation problem becomes that of **decomposing** \mathbf{X} into **blocks** of the form **(low-rank matrix) \otimes vector**:

$$\mathbf{X} = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{w}_r, \quad \forall r, \text{rank } \mathbf{H}_r \leq L_r.$$

As before, it turns out that $\{(\mathbf{H}_r, \mathbf{w}_r)\}_{r=1}^R$ are **(essentially) unique** under mild constraints.¹

1: De Lathauwer, 2011

Example #3: Multilinear PCA

We're given $M \times M$ face images of N_p people:

$$\mathbf{X} \in \mathbb{R}^{N \times N_p}, \quad \text{with } N = M^2.$$

Dimensionality reduction and feature extraction can be achieved by [PCA](#) (via truncated SVD of \mathbf{X} , after centering):

$$\mathbf{X} \approx \mathbf{U}_R \mathbf{\Sigma}_R \mathbf{V}_R^\top = \sum_{r=1}^R \sigma_r(\mathbf{U}_R)_{:r} (\mathbf{V}_R)_{:r}^\top, \quad x_{ip} \approx \sum_{r=1}^R \sigma_r(\mathbf{U}_R)_{ir} (\mathbf{V}_R)_{pr},$$

where

- $\mathbf{U}_R \in \mathbb{R}^{N \times R}$ contains the R first left singular vectors
- $\mathbf{V}_R \in \mathbb{R}^{N_p \times R}$ contains the R first right singular vectors
- $\mathbf{\Sigma}_R = \text{Diag}(\sigma_1, \dots, \sigma_R)$ contains the R first singular values

Example #3: Multilinear PCA

We're given $M \times M$ face images of N_p people:

$$\mathbf{X} \in \mathbb{R}^{N \times N_p}, \quad \text{with } N = M^2.$$

Dimensionality reduction and feature extraction can be achieved by [PCA](#) (via truncated SVD of \mathbf{X} , after centering):

$$\mathbf{X} \approx \mathbf{U}_R \mathbf{\Sigma}_R \mathbf{V}_R^\top = \sum_{r=1}^R \sigma_r(\mathbf{U}_R)_{:r} (\mathbf{V}_R)_{:r}^\top, \quad x_{ip} \approx \sum_{r=1}^R \sigma_r(\mathbf{U}_R)_{ir} (\mathbf{V}_R)_{pr},$$

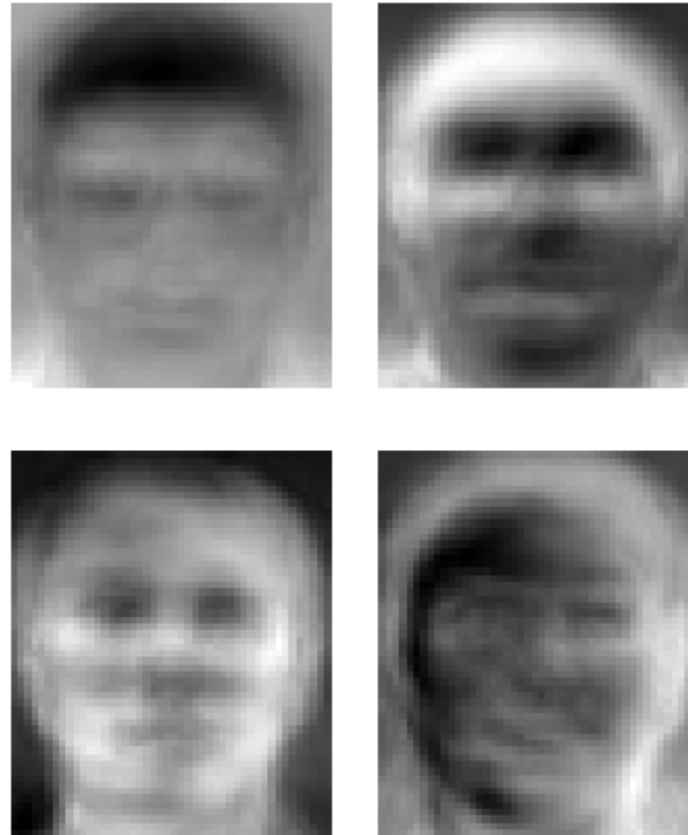
where

- $\mathbf{U}_R \in \mathbb{R}^{N \times R}$ contains the R first left singular vectors
- $\mathbf{V}_R \in \mathbb{R}^{N_p \times R}$ contains the R first right singular vectors
- $\mathbf{\Sigma}_R = \text{Diag}(\sigma_1, \dots, \sigma_R)$ contains the R first singular values

Hence, $(\mathbf{\Sigma}_R \mathbf{V}_R^\top)_{:p} \in \mathbb{R}^R$ contains the coordinates of the p th person's image w.r.t. the subspace basis \mathbf{U}_R .

Eigenfaces

This is the so-called **eigenfaces**¹ approach to dimension reduction and feature extraction.



A few eigenfaces (cols of \mathbf{U}_R) from the ORL database (source: Wikipedia).

1: Sirovich & Kirby, 1987

Extension to multiple diversities

Consider next N -dim face images of N_p people under N_i illumination conditions.

PCA of $\mathbf{X} \in \mathbb{R}^{N \times N_i N_p}$ is still possible, but two diversities “get entangled.”

Extension to multiple diversities

Consider next N -dim face images of N_p people under N_i illumination conditions.

PCA of $\mathbf{X} \in \mathbb{R}^{N \times N_i N_p}$ is still possible, but two diversities “get entangled.”

Idea: Disentangle them via a tensor structured basis

$$\mathbf{X} \approx \mathbf{U} \mathbf{S} (\mathbf{U}^{(i)} \boxtimes \mathbf{U}^{(p)})^T.$$

Def: Kronecker product $\mathbf{A} \boxtimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \dots \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$

Extension to multiple diversities

Consider next N -dim face images of N_p people under N_i illumination conditions.

PCA of $\mathbf{X} \in \mathbb{R}^{N \times N_i N_p}$ is still possible, but two diversities “get entangled.”

Idea: Disentangle them via a tensor structured basis

$$\mathbf{X} \approx \mathbf{U} \mathbf{S} (\mathbf{U}^{(i)} \boxtimes \mathbf{U}^{(p)})^\top.$$

Def: Kronecker product $\mathbf{A} \boxtimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \dots \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$

In tensor notation: the data $\mathbf{X} \in \mathbb{R}^{N \times N_p \times N_i}$ is modeled as

$$\mathbf{X} \approx (\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}) \cdot \mathbf{S} := \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} s_{r_1 r_2 r_3} (\mathbf{U})_{:r_1} \otimes (\mathbf{U}^{(p)})_{:r_2} \otimes (\mathbf{U}^{(i)})_{:r_3}$$

Multilinear PCA

We can see $(\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}) \cdot \mathbf{S}$ as a **contraction** of \mathbf{S} with $\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}$:

$$x_{npi} \approx \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (\mathbf{U})_{nr_1} (\mathbf{U}^{(p)})_{pr_2} (\mathbf{U}^{(i)})_{ir_3}.$$

Multilinear PCA

We can see $(\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}) \cdot \mathbf{S}$ as a **contraction** of \mathbf{S} with $\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}$:

$$x_{npi} \approx \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (\mathbf{U})_{nr_1} (\mathbf{U}^{(p)})_{pr_2} (\mathbf{U}^{(i)})_{ir_3}.$$

Hence, the partial contraction $(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S} \in \mathbb{R}^{N \times R_2 \times R_3}$ given by

$$[(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S}]_{nr_2 r_3} = \sum_{r_1} s_{r_1 r_2 r_3} (\mathbf{U})_{nr_1}$$

yields a tensor having $R_2 R_3$ images which represent R_2 “people patterns” $\times R_3$ “illumination patterns.”

We can see $(\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}) \cdot \mathbf{S}$ as a **contraction** of \mathbf{S} with $\mathbf{U}, \mathbf{U}^{(p)}, \mathbf{U}^{(i)}$:

$$x_{npi} \approx \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (\mathbf{U})_{nr_1} (\mathbf{U}^{(p)})_{pr_2} (\mathbf{U}^{(i)})_{ir_3}.$$

Hence, the partial contraction $(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S} \in \mathbb{R}^{N \times R_2 \times R_3}$ given by

$$[(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S}]_{nr_2 r_3} = \sum_{r_1} s_{r_1 r_2 r_3} (\mathbf{U})_{nr_1}$$

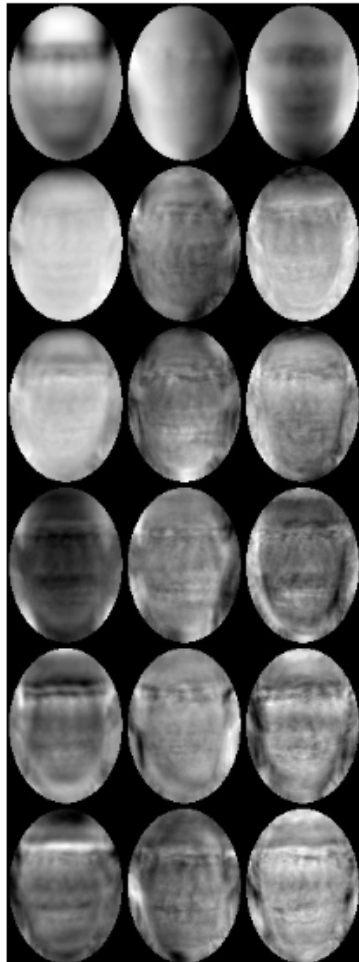
yields a tensor having $R_2 R_3$ images which represent R_2 “people patterns” $\times R_3$ “illumination patterns.”

Further contraction with the p th row of $\mathbf{U}^{(p)}$ and the i th row of $\mathbf{U}^{(i)}$ gives the approximate image of person p under illumination condition i :

$$(\mathbf{X})_{:pi} \approx \sum_{r_2 r_3} \underbrace{\left(\sum_{r_1} s_{r_1 r_2 r_3} (\mathbf{U})_{:r_1} \right)}_{\text{image pattern } (r_2, r_3)} (\mathbf{U}^{(p)})_{pr_2} (\mathbf{U}^{(i)})_{ir_3}.$$

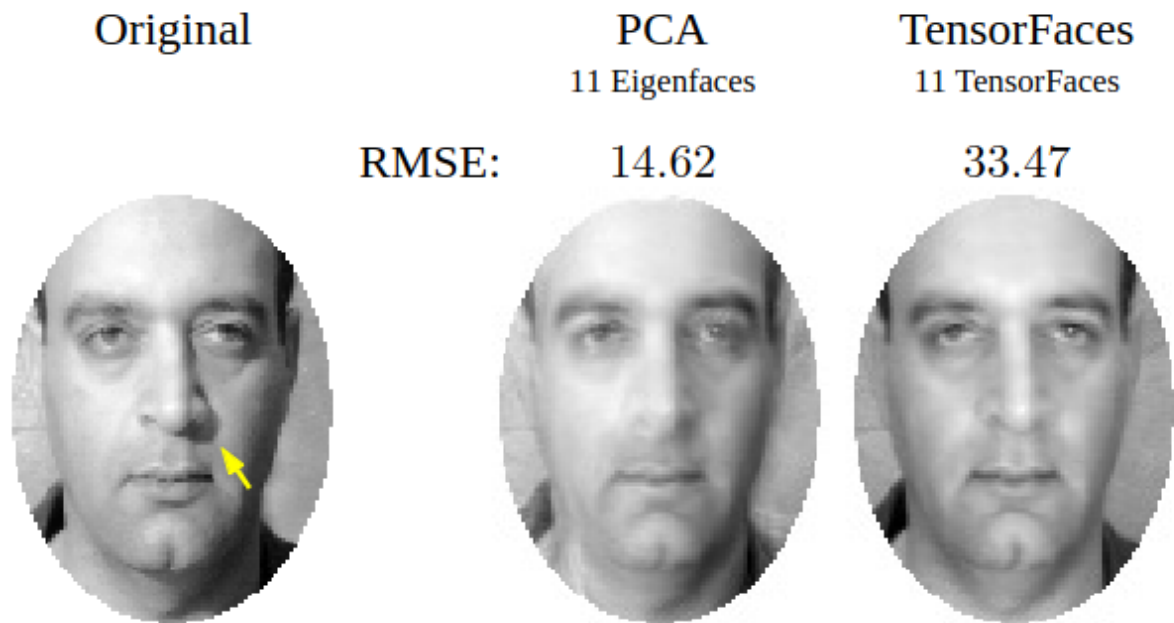
This idea forms the basis of the **TensorFaces** approach¹.

people ↓ illuminations →



A few “eigenmodes.”¹

It allows in particular a “strategic” dim. reduction w.r.t. some chosen diversities (e.g., illumination only, via a small R_3).

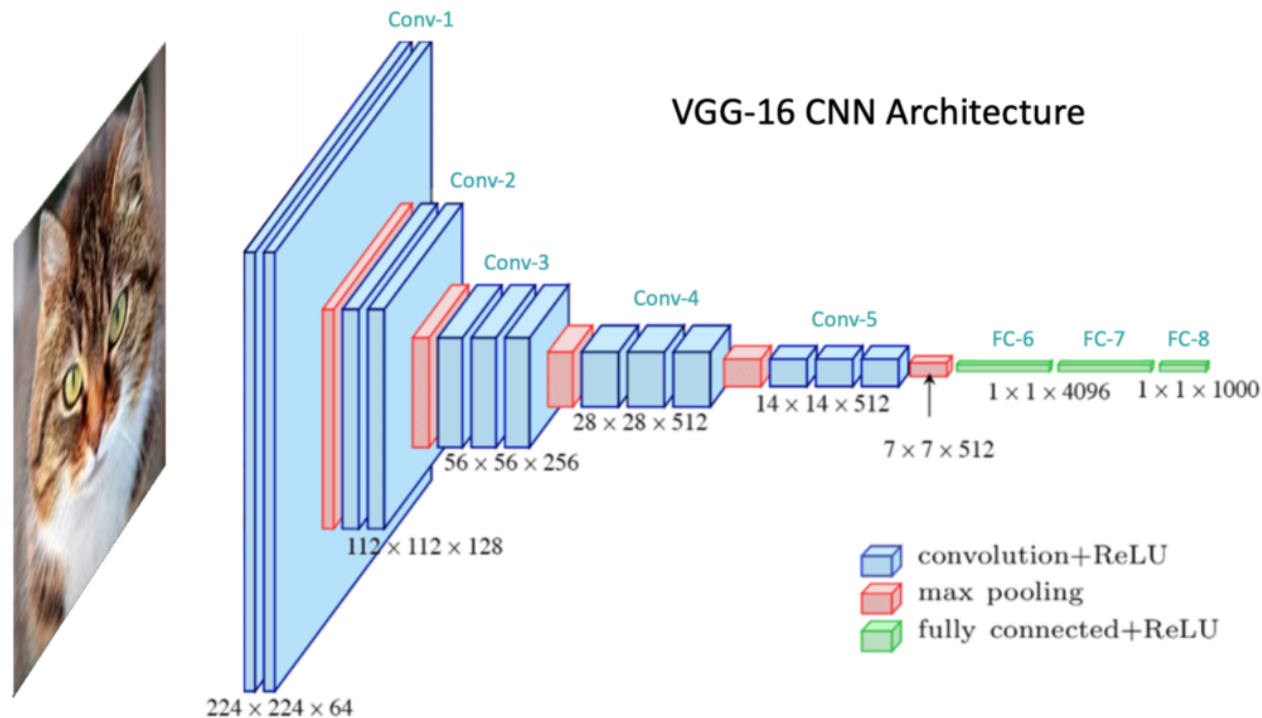


*Dimension reduction:*¹ PCA vs TensorFaces with $(R_1, R_2, R_3) = (176, 11, 1)$. (Original size: $7943 \times 11 \times 16$.)

Example #4: Deep learning (of course!)

Numerous deep learning models are parameterized by tensors.

Typical example: CNNs parameterized by **convolution kernels**.



A typical CNN architecture (source: learnopencv.com).

Kernel of conv2D layer with $H \times W$ filters, I input channels and O output channels:

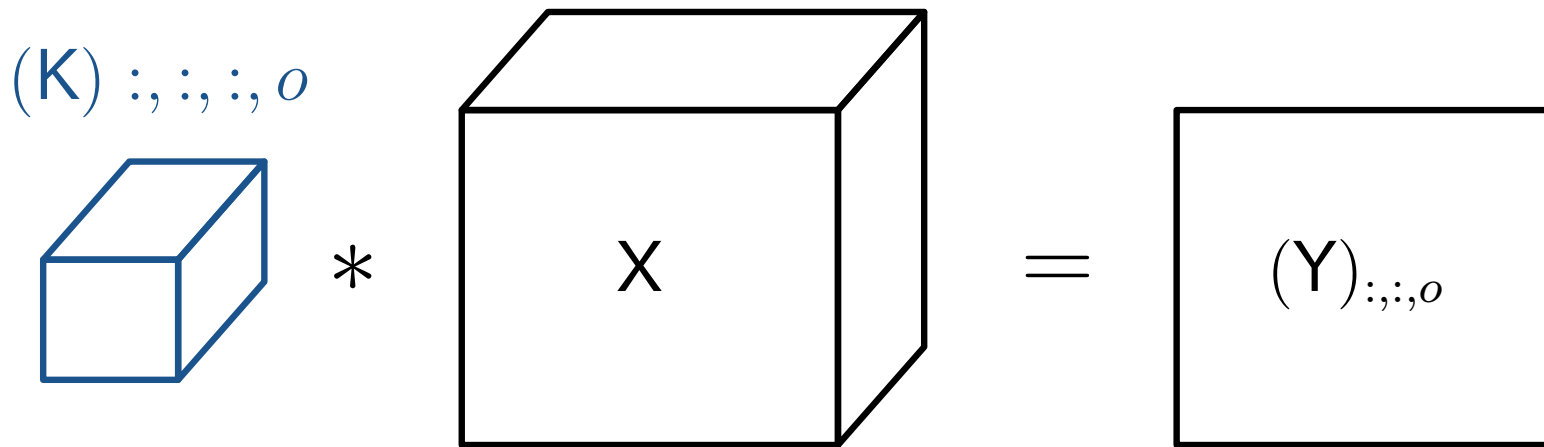
$$\mathbf{K} \in \mathbb{R}^{H \times W \times I \times O}.$$

2D convolution in CNNs

conv2D maps $\mathbf{X} \in \mathbb{R}^{M \times N \times I}$ into $\mathbf{Y} \in \mathbb{R}^{M \times N \times O}$ according to

$$y_{mno} = \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \sum_{i=1}^I (\mathbf{K})_{hwo} x_{m+h, n+w, i}.$$

Per output channel o :



Complexity: $\mathcal{O}(MNOHWI)$.

Most expensive stage of inference: performing the required convolutions.

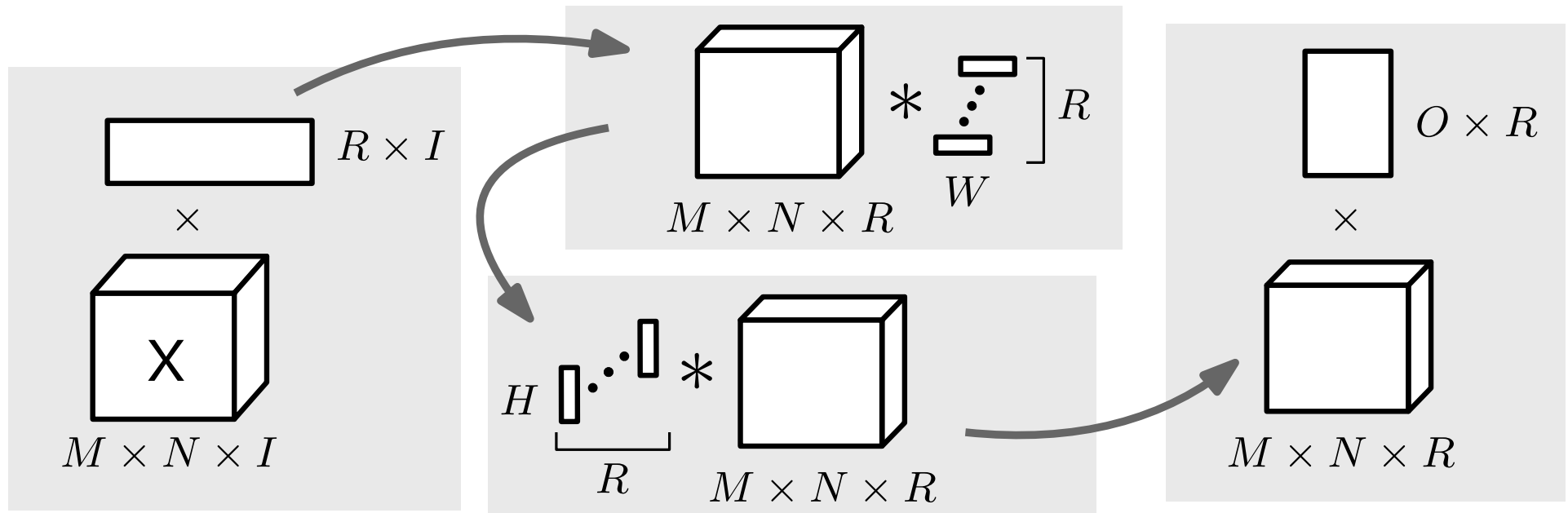
Decomposing convolution kernels

Idea: decompose K into separable terms using $R(H + W + I + O)$ params,¹

$$(K)_{hwio} = \sum_{r=1}^R a_{hr} b_{wr} c_{ir} d_{or},$$

$$y_{mno} = \sum_{r=1}^R d_{or} \sum_{h=0}^{H-1} a_{hr} \sum_{w=0}^{W-1} b_{wr} \sum_{i=1}^I c_{ir} x_{m+h,w+n,i}.$$

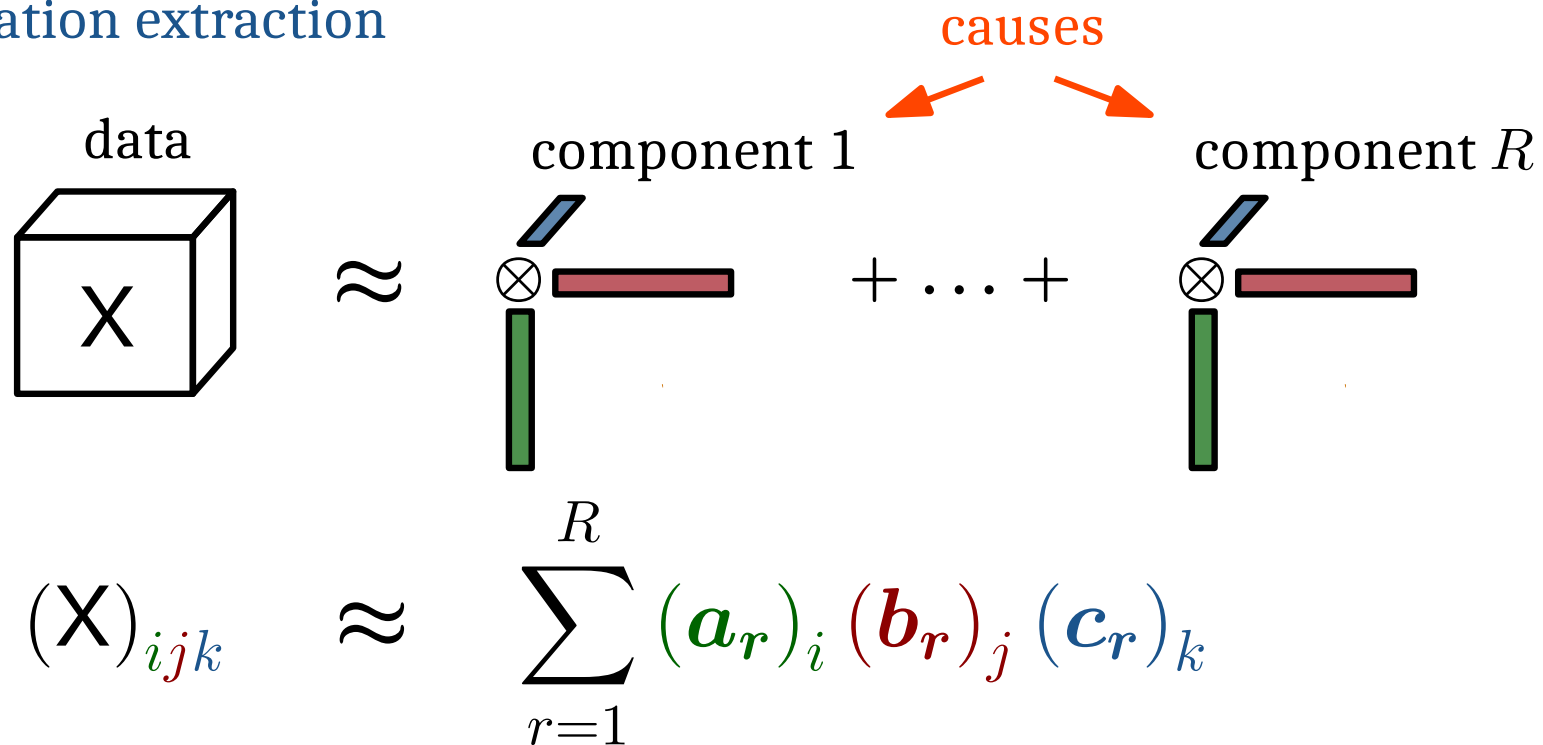
Complexity: $\mathcal{O}(MNR(H + W + I + O))$



Summary : two main motivations (1/2)

Tensor models generally serve two main purposes:

1. Information extraction



data = observations or functions thereof, R “small”, “simple” components

Ex: exploratory data analysis,¹ source separation, latent variable model estimation, ...

1: Hong & al., 2020

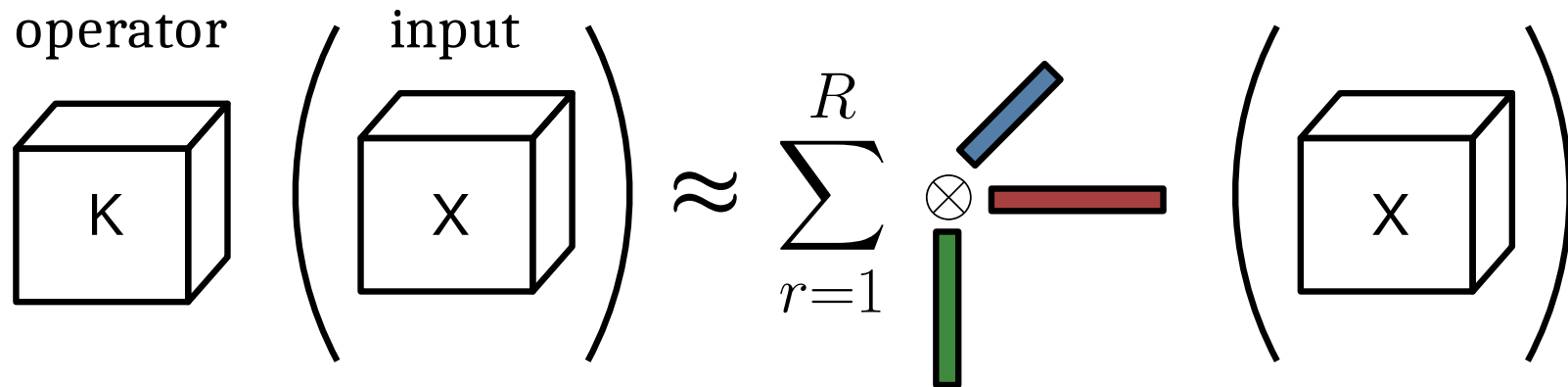
Summary : two main motivations (2/2)

2. Complexity reduction

- Dimensionality reduction with tensor structured basis

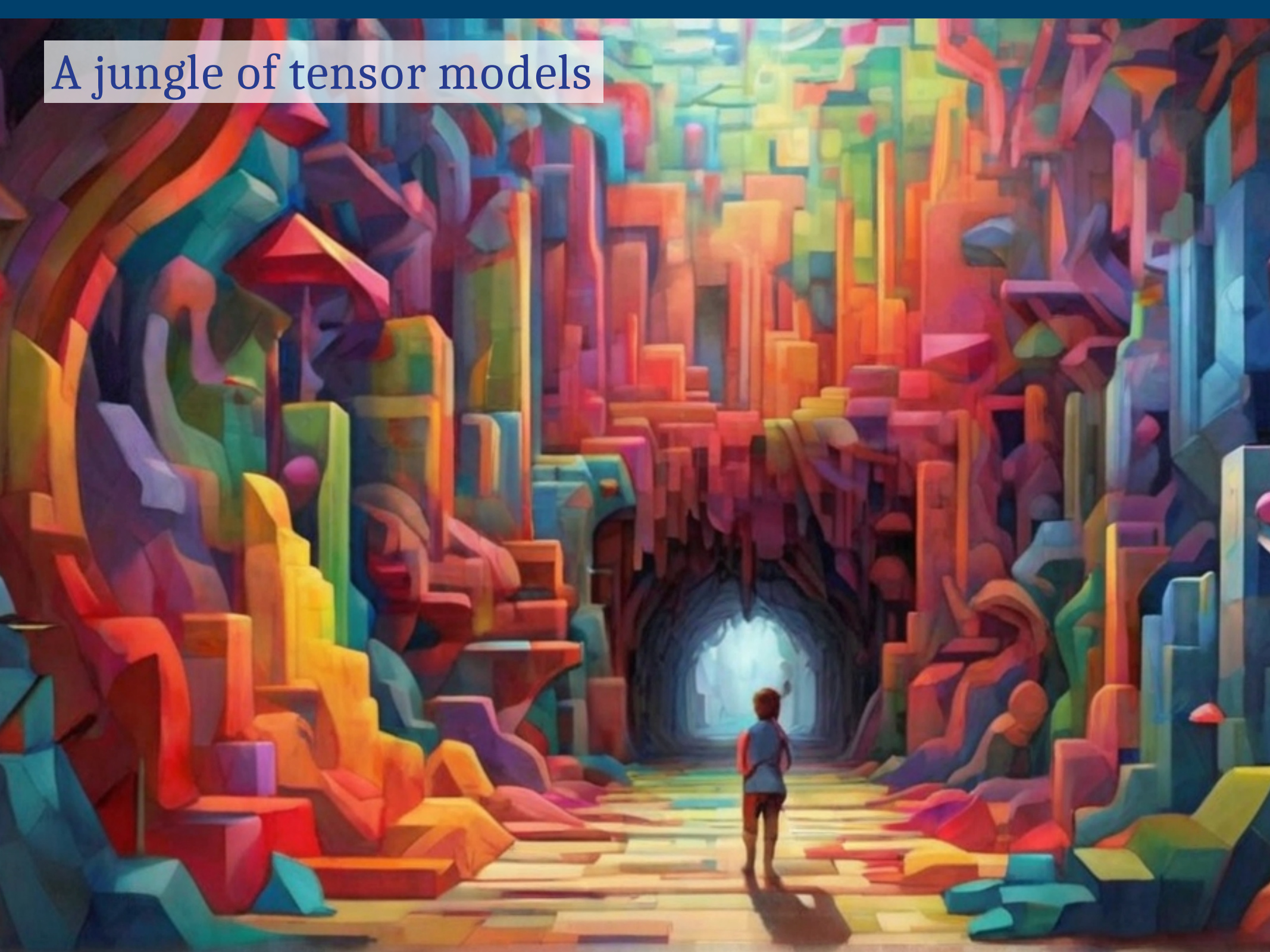
$$\mathbf{X} \approx (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathbf{S}$$

- Operator \approx operator with “low-rank structure”



Ex: nonlinear system modelling,^{1,2} deep learning, multivariate density approximation,³ ...

A jungle of tensor models

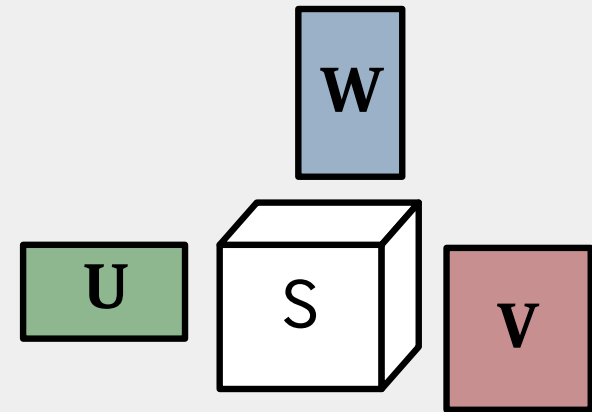


Let's start with some useful notation & definitions.

Multilinear transformation

Def: Given $\mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$ and $\mathbf{U} \in \mathbb{R}^{N_1 \times R_1}$, $\mathbf{V} \in \mathbb{R}^{N_2 \times R_2}$ and $\mathbf{W} \in \mathbb{R}^{N_3 \times R_3}$, the **multilinear transformation**

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

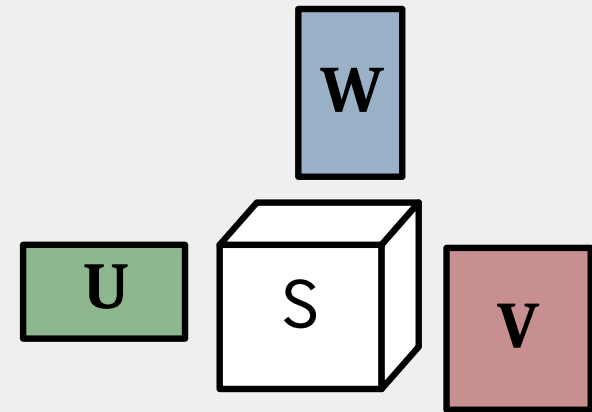


is defined as
$$[(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}]_{n_1 n_2 n_3} = \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (\mathbf{U})_{n_1 r_1} (\mathbf{V})_{n_2 r_2} (\mathbf{W})_{n_3 r_3}.$$

Multilinear transformation

Def: Given $\mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$ and $\mathbf{U} \in \mathbb{R}^{N_1 \times R_1}$, $\mathbf{V} \in \mathbb{R}^{N_2 \times R_2}$ and $\mathbf{W} \in \mathbb{R}^{N_3 \times R_3}$, the **multilinear transformation**

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$



is defined as
$$[(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}]_{n_1 n_2 n_3} = \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (\mathbf{U})_{n_1 r_1} (\mathbf{V})_{n_2 r_2} (\mathbf{W})_{n_3 r_3}.$$

Notational conventions:

(i) replacing any matrix by a dot \cdot means no transformation in that mode

$$[(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S}]_{n_1 r_2 r_3} = \sum_{r_1} s_{r_1 r_2 r_3} (\mathbf{U})_{n_1 r_1}$$

(ii) the transpose is omitted when transformation by vectors are performed

$$(u, v, w) \cdot \mathbf{S} = \sum_{r_1 r_2 r_3} s_{r_1 r_2 r_3} (u)_{r_1} (v)_{r_2} (w)_{r_3}.$$

Multilinear transformation: properties

(i) Expansion in terms of a tensor basis:

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} s_{r_1 r_2 r_3} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3}$$

(ii) Action on rank-1 tensors:

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} = (\mathbf{U}\mathbf{a}) \otimes (\mathbf{V}\mathbf{b}) \otimes (\mathbf{W}\mathbf{c})$$

Hence, by linearity,

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \left(\sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r \right) = \sum_{r=1}^R (\mathbf{U}\mathbf{a}_r) \otimes (\mathbf{V}\mathbf{b}_r) \otimes (\mathbf{W}\mathbf{c}_r)$$

(iii) Composition:

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot [(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \cdot \mathbf{S}] = (\mathbf{UX}, \mathbf{VY}, \mathbf{WZ}) \cdot \mathbf{S}$$

$$(\mathbf{U}, \cdot, \cdot) \cdot [(\cdot, \mathbf{V}, \cdot) \cdot \mathbf{S}] = (\cdot, \mathbf{V}, \cdot) \cdot [(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{S}]$$

Scalar product and Euclidean norm

Def: Scalar product

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_{ijk} x_{ijk} y_{ijk}$$

In particular:
$$\begin{cases} \langle \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}, \mathbf{Y} \rangle = (\mathbf{a}, \mathbf{b}, \mathbf{c}) \cdot \mathbf{Y} \\ \langle \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}, \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} \rangle = \langle \mathbf{a}, \mathbf{u} \rangle \langle \mathbf{b}, \mathbf{v} \rangle \langle \mathbf{c}, \mathbf{w} \rangle \end{cases}$$

Def: The Frobenius norm of \mathbf{W} is defined as $\|\mathbf{W}\|_F := \sqrt{\langle \mathbf{W}, \mathbf{W} \rangle} = \sqrt{\sum_{ijk} w_{ijk}^2}$

Natural (& useful) isomorphisms

Idea: “view” tensors through the lens of certain isomorphisms with matrix and vector spaces. Useful for both **analysis & computation**.

(Two finite-dim vector spaces are **isomorphic** (\simeq) iff they have the same dim.)

Natural (& useful) isomorphisms

Idea: “view” tensors through the lens of certain isomorphisms with matrix and vector spaces. Useful for both **analysis & computation**.

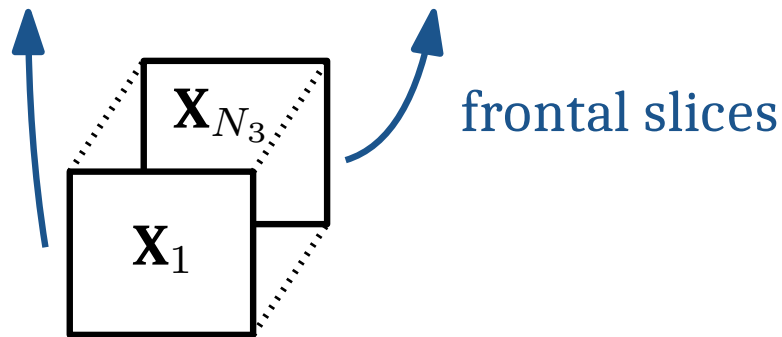
(Two finite-dim vector spaces are **isomorphic** (\simeq) iff they have the same dim.)

In particular, the isomorphism

$$\mathbb{R}^{N_1 \times N_2 \times N_3} \simeq \mathbb{R}^{N_1 \times N_3 N_2}.$$

leads to the notion of a mode-1 **unfolding**

$$\mathbf{X} \mapsto \mathbf{X}_{(1)} := (\mathbf{X}_1 \quad \mathbf{X}_2 \quad \dots \quad \mathbf{X}_{N_3}) \in \mathbb{R}^{N_1 \times N_3 N_2}$$



Concretely, x_{ijk} is the element $(i, n(j, k))$ of $\mathbf{X}_{(1)}$, with $n(j, k) := (k - 1)N_2 + j$.

Analogous definitions can be given to $\mathbf{X}_{(2)} \in \mathbb{R}^{N_2 \times N_3 N_1}$ and $\mathbf{X}_{(3)} \in \mathbb{R}^{N_3 \times N_2 N_1}$.

Natural (& useful) isomorphisms (cont'd)

Ex: If $X = a \otimes b \otimes c$, then:

$$\mathbf{X}_{(1)} = a (c \boxtimes b)^\top, \quad \mathbf{X}_{(2)} = b (c \boxtimes a)^\top, \quad \mathbf{X}_{(3)} = c (b \boxtimes a)^\top.$$

Natural (& useful) isomorphisms (cont'd)

Ex: If $X = a \otimes b \otimes c$, then:

$$\mathbf{X}_{(1)} = a (c \boxtimes b)^\top, \quad \mathbf{X}_{(2)} = b (c \boxtimes a)^\top, \quad \mathbf{X}_{(3)} = c (b \boxtimes a)^\top.$$

Ex: If $X = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}$, then:

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top, \quad \mathbf{X}_{(2)} = \mathbf{V} \mathbf{S}_{(2)} (\mathbf{W} \boxtimes \mathbf{U})^\top, \quad \mathbf{X}_{(3)} = \mathbf{W} \mathbf{S}_{(3)} (\mathbf{V} \boxtimes \mathbf{U})^\top.$$

Natural (& useful) isomorphisms (cont'd)

Ex: If $X = a \otimes b \otimes c$, then:

$$\mathbf{X}_{(1)} = a (c \boxtimes b)^\top, \quad \mathbf{X}_{(2)} = b (c \boxtimes a)^\top, \quad \mathbf{X}_{(3)} = c (b \boxtimes a)^\top.$$

Ex: If $X = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}$, then:

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top, \quad \mathbf{X}_{(2)} = \mathbf{V} \mathbf{S}_{(2)} (\mathbf{W} \boxtimes \mathbf{U})^\top, \quad \mathbf{X}_{(3)} = \mathbf{W} \mathbf{S}_{(3)} (\mathbf{V} \boxtimes \mathbf{U})^\top.$$

Rmk: In practice, $X = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}$ is computed by a sequence of matrix-matrix products and unfoldings/foldings.

Natural (& useful) isomorphisms (cont'd)

Ex: If $\mathbf{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$, then:

$$\mathbf{X}_{(1)} = \mathbf{a} (\mathbf{c} \boxtimes \mathbf{b})^\top, \quad \mathbf{X}_{(2)} = \mathbf{b} (\mathbf{c} \boxtimes \mathbf{a})^\top, \quad \mathbf{X}_{(3)} = \mathbf{c} (\mathbf{b} \boxtimes \mathbf{a})^\top.$$

Ex: If $\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}$, then:

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top, \quad \mathbf{X}_{(2)} = \mathbf{V} \mathbf{S}_{(2)} (\mathbf{W} \boxtimes \mathbf{U})^\top, \quad \mathbf{X}_{(3)} = \mathbf{W} \mathbf{S}_{(3)} (\mathbf{V} \boxtimes \mathbf{U})^\top.$$

Rmk: In practice, $\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}$ is computed by a sequence of matrix-matrix products and unfoldings/foldings.

Another commonly used isomorphism is the **vectorization** operation

$$\mathbb{R}^{N_1 \times N_2 \times N_3} \ni \mathbf{X} \mapsto \mathbf{x} = \text{vec } \mathbf{X} \in \mathbb{R}^{N_3 N_2 N_1},$$

such that \mathbf{X}_{ijk} is “sent” to $x_{n(i,j,k)}$ with $n(i, j, k) := N_2 N_1 (k - 1) + N_1 (j - 1) + i$.

Ex: By vectorization,

$$\text{vec} (\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}) = \mathbf{c} \boxtimes \mathbf{b} \boxtimes \mathbf{a},$$

$$\text{vec} [(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}] = (\mathbf{W} \boxtimes \mathbf{V} \boxtimes \mathbf{U}) \text{vec } \mathbf{S}.$$

...and now let's finally discuss low-rank tensor models.

What is the rank of a tensor?

For a matrix $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$,

$$\begin{aligned} \text{rank } \mathbf{X} &= \dim \text{colspan } \mathbf{X} = \dim \left\{ \mathbf{X} \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_2} \right\} \\ &= \dim \text{rowspan } \mathbf{X} = \dim \left\{ \mathbf{X}^\top \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_1} \right\} \\ &= \min \left\{ R : \mathbf{X} = \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^\top \right\} \end{aligned}$$

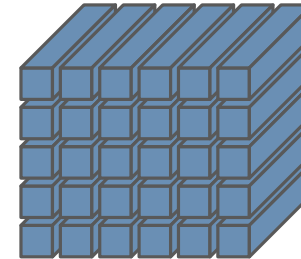
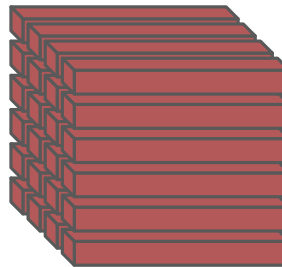
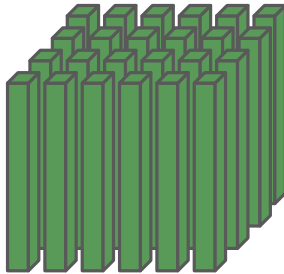
What is the rank of a tensor?

For a matrix $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$,

$$\begin{aligned} \text{rank } \mathbf{X} &= \dim \text{colspan } \mathbf{X} = \dim \{ \mathbf{X} \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_2} \} \\ &= \dim \text{rowspan } \mathbf{X} = \dim \{ \mathbf{X}^\top \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_1} \} \\ &= \min \left\{ R : \mathbf{X} = \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^\top \right\} \end{aligned}$$

This suggests (at least) two definitions for the rank of $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$:

1. Dimension(s) of subspaces spanned by the **fibers** of each mode:



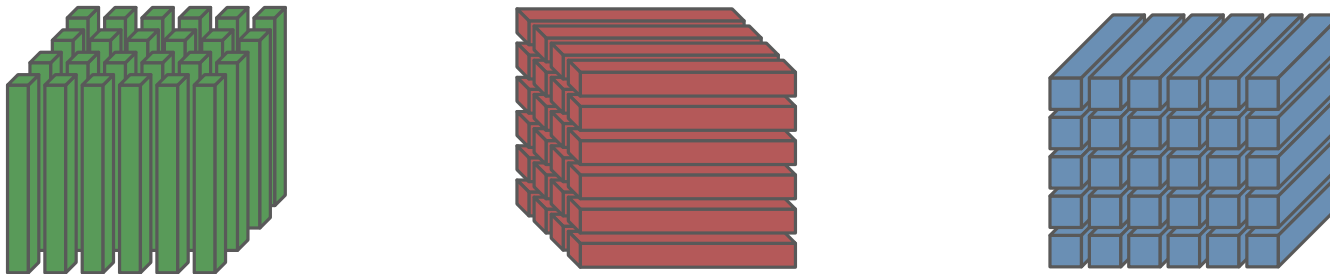
What is the rank of a tensor?

For a matrix $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$,

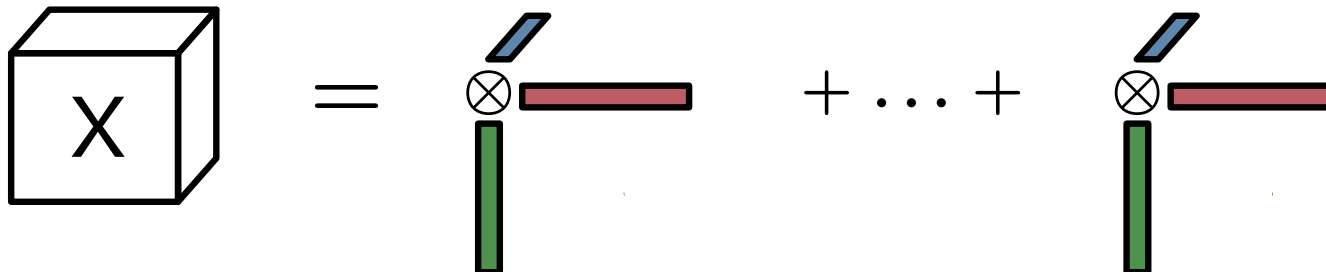
$$\begin{aligned} \text{rank } \mathbf{X} &= \dim \text{colspan } \mathbf{X} = \dim \{ \mathbf{X} \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_2} \} \\ &= \dim \text{rowspan } \mathbf{X} = \dim \{ \mathbf{X}^\top \mathbf{v} : \mathbf{v} \in \mathbb{R}^{N_1} \} \\ &= \min \left\{ R : \mathbf{X} = \sum_{r=1}^R \mathbf{a}_r \mathbf{b}_r^\top \right\} \end{aligned}$$

This suggests (at least) two definitions for the rank of $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$:

1. Dimension(s) of subspaces spanned by the **fibers** of each mode:



2. Minimal number of rank-1 tensors needed to produce \mathbf{X} :



1. The multilinear rank

Recall that: fibers of \mathbf{X} = columns of each $\mathbf{X}_{(i)}$.

1. The multilinear rank

Recall that: fibers of \mathbf{X} = columns of each $\mathbf{X}_{(i)}$.

The first notion of rank can be formalized as

$$\begin{aligned}\text{rank}_1 \mathbf{X} &:= \dim \text{span} \{ (\cdot, \mathbf{v}, \mathbf{w}) \cdot \mathbf{X} : \mathbf{v} \in \mathbb{R}^{N_2}, \mathbf{w} \in \mathbb{R}^{N_3} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(1)}(\mathbf{w} \boxtimes \mathbf{v}) \} = \text{rank } \mathbf{X}_{(1)}\end{aligned}$$

$$\begin{aligned}\text{rank}_2 \mathbf{X} &:= \dim \text{span} \{ (\mathbf{u}, \cdot, \mathbf{w}) \cdot \mathbf{X} : \mathbf{u} \in \mathbb{R}^{N_1}, \mathbf{w} \in \mathbb{R}^{N_3} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(2)}(\mathbf{w} \boxtimes \mathbf{u}) \} = \text{rank } \mathbf{X}_{(2)}\end{aligned}$$

$$\begin{aligned}\text{rank}_3 \mathbf{X} &:= \dim \text{span} \{ (\mathbf{u}, \mathbf{v}, \cdot) \cdot \mathbf{X} : \mathbf{u} \in \mathbb{R}^{N_1}, \mathbf{v} \in \mathbb{R}^{N_2} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(3)}(\mathbf{v} \boxtimes \mathbf{u}) \} = \text{rank } \mathbf{X}_{(3)}\end{aligned}$$

1. The multilinear rank

Recall that: fibers of \mathbf{X} = columns of each $\mathbf{X}_{(i)}$.

The first notion of rank can be formalized as

$$\begin{aligned} \text{rank}_1 \mathbf{X} &:= \dim \text{span} \{ (\cdot, \mathbf{v}, \mathbf{w}) \cdot \mathbf{X} : \mathbf{v} \in \mathbb{R}^{N_2}, \mathbf{w} \in \mathbb{R}^{N_3} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(1)}(\mathbf{w} \boxtimes \mathbf{v}) \} = \text{rank } \mathbf{X}_{(1)} \end{aligned}$$

$$\begin{aligned} \text{rank}_2 \mathbf{X} &:= \dim \text{span} \{ (\mathbf{u}, \cdot, \mathbf{w}) \cdot \mathbf{X} : \mathbf{u} \in \mathbb{R}^{N_1}, \mathbf{w} \in \mathbb{R}^{N_3} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(2)}(\mathbf{w} \boxtimes \mathbf{u}) \} = \text{rank } \mathbf{X}_{(2)} \end{aligned}$$

$$\begin{aligned} \text{rank}_3 \mathbf{X} &:= \dim \text{span} \{ (\mathbf{u}, \mathbf{v}, \cdot) \cdot \mathbf{X} : \mathbf{u} \in \mathbb{R}^{N_1}, \mathbf{v} \in \mathbb{R}^{N_2} \} \\ &= \dim \text{span} \{ \mathbf{X}_{(3)}(\mathbf{v} \boxtimes \mathbf{u}) \} = \text{rank } \mathbf{X}_{(3)} \end{aligned}$$

Def:¹ Multilinear rank: $\text{mrnk } \mathbf{X} := (\text{rank}_1 \mathbf{X}, \text{rank}_2 \mathbf{X}, \text{rank}_3 \mathbf{X})$.

1: Hitchcock, 1928

Properties of mrank

(i) Unlike the matrix case, mrank components are **generally distinct**.

Ex: With

$$\mathbf{X} = \left(\begin{array}{cc|cc} & \mathbf{X}_1 & & \mathbf{X}_2 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{array} \right) \in \mathbb{R}^{2 \times 2 \times 2},$$

$$\mathbf{X}_{(1)} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad \mathbf{X}_{(2)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad \mathbf{X}_{(3)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

Properties of mrank

(i) Unlike the matrix case, mrank components are **generally distinct**.

Ex: With

$$\mathbf{X} = \left(\begin{array}{cc|cc} & \mathbf{X}_1 & & \mathbf{X}_2 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{array} \right) \in \mathbb{R}^{2 \times 2 \times 2},$$

$$\mathbf{X}_{(1)} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad \mathbf{X}_{(2)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad \mathbf{X}_{(3)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

(ii) The mrank is a tensor property:

Prop: mrank \mathbf{X} is invariant w.r.t. a change of basis

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{X}, \quad \mathbf{U} \in \text{GL}(N_1), \mathbf{V} \in \text{GL}(N_2), \mathbf{W} \in \text{GL}(N_3).$$

Proof: For any nonsingular $\mathbf{U} \in \mathbb{R}^{N_1 \times N_1}$ (and similarly for the other modes),

$$\text{rank}_1[(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{X}] = \text{rank}[\mathbf{U}\mathbf{X}_{(1)}] = \text{rank } \mathbf{X}_{(1)} = \text{rank}_1 \mathbf{X}. \quad \blacksquare$$

Properties of mrank

(i) Unlike the matrix case, mrank components are **generally distinct**.

Ex: With

$$\mathbf{X} = \begin{pmatrix} & \mathbf{X}_1 & & \mathbf{X}_2 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{2 \times 2 \times 2},$$

$$\mathbf{X}_{(1)} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad \mathbf{X}_{(2)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad \mathbf{X}_{(3)} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

(ii) The mrank is a tensor property:

Prop: mrank \mathbf{X} is invariant w.r.t. a change of basis

$$(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{X}, \quad \mathbf{U} \in \text{GL}(N_1), \mathbf{V} \in \text{GL}(N_2), \mathbf{W} \in \text{GL}(N_3).$$

Proof: For any nonsingular $\mathbf{U} \in \mathbb{R}^{N_1 \times N_1}$ (and similarly for the other modes),

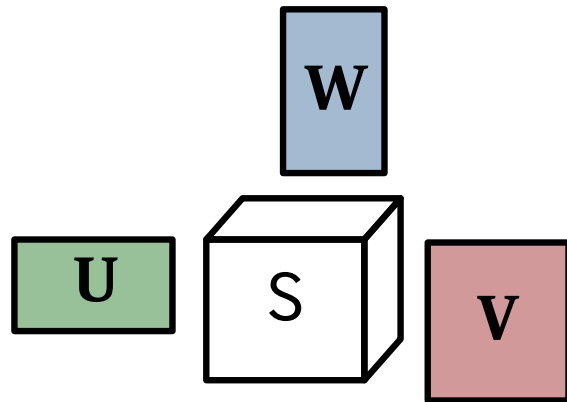
$$\text{rank}_1[(\mathbf{U}, \cdot, \cdot) \cdot \mathbf{X}] = \text{rank}[\mathbf{U}\mathbf{X}_{(1)}] = \text{rank} \mathbf{X}_{(1)} = \text{rank}_1 \mathbf{X}. \quad \blacksquare$$

(iii) $\text{mrk} \mathbf{X} \leq (\min \{N_1, N_2 N_3\}, \min \{N_2, N_1 N_3\}, \min \{N_3, N_1 N_2\})$

The Tucker model

The Tucker model¹ has the form

$$\mathbb{R}^{N_1 \times N_2 \times N_3} \ni \mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$$



with:

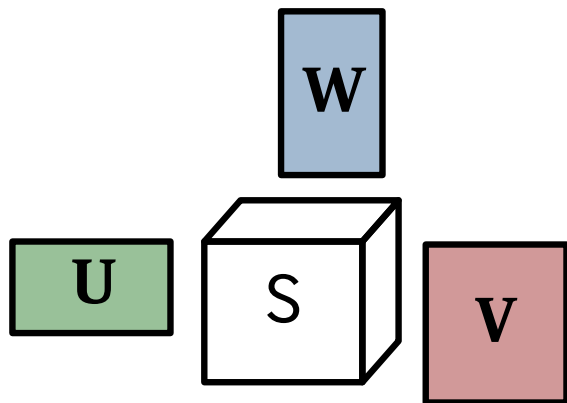
- $\mathbf{U} \in \mathbb{R}^{N_1 \times R_1}, R_1 \leq N_1$
- $\mathbf{V} \in \mathbb{R}^{N_2 \times R_2}, R_2 \leq N_2$
- $\mathbf{W} \in \mathbb{R}^{N_3 \times R_3}, R_3 \leq N_3$

1: Tucker, 1966

The Tucker model

The Tucker model¹ has the form

$$\mathbb{R}^{N_1 \times N_2 \times N_3} \ni \mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$$



with:

- $\mathbf{U} \in \mathbb{R}^{N_1 \times R_1}, R_1 \leq N_1$
- $\mathbf{V} \in \mathbb{R}^{N_2 \times R_2}, R_2 \leq N_2$
- $\mathbf{W} \in \mathbb{R}^{N_3 \times R_3}, R_3 \leq N_3$

Recall that

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top$$

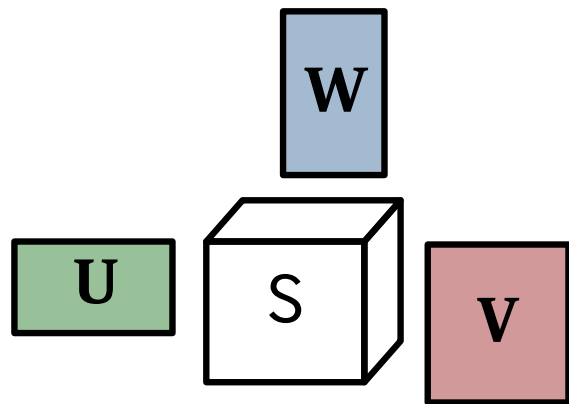
$$\mathbf{X}_{(2)} = \mathbf{V} \mathbf{S}_{(2)} (\mathbf{W} \boxtimes \mathbf{U})^\top$$

$$\mathbf{X}_{(3)} = \mathbf{W} \mathbf{S}_{(3)} (\mathbf{V} \boxtimes \mathbf{U})^\top$$

The Tucker model

The Tucker model¹ has the form

$$\mathbb{R}^{N_1 \times N_2 \times N_3} \ni \mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3}$$



with:

- $\mathbf{U} \in \mathbb{R}^{N_1 \times R_1}, R_1 \leq N_1$
- $\mathbf{V} \in \mathbb{R}^{N_2 \times R_2}, R_2 \leq N_2$
- $\mathbf{W} \in \mathbb{R}^{N_3 \times R_3}, R_3 \leq N_3$

Recall that

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top$$

$$\mathbf{X}_{(2)} = \mathbf{V} \mathbf{S}_{(2)} (\mathbf{W} \boxtimes \mathbf{U})^\top$$

$$\mathbf{X}_{(3)} = \mathbf{W} \mathbf{S}_{(3)} (\mathbf{V} \boxtimes \mathbf{U})^\top$$

Clearly, $\text{mrnk } \mathbf{X} \leq (R_1, R_2, R_3)$ component-wise.

Hence, useful for a **low-mrank** representation of a tensor (**multilinear PCA**).

1: Tucker, 1966

Orthogonal Tucker model

The Tucker model is a **subspace representation** of a tensor w.r.t. the **tensor basis**¹

$$\{\mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3} : r_i \in [R_i], 1 \leq i \leq 3\},$$

since

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} s_{r_1 r_2 r_3} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3}.$$

Orthogonal Tucker model

The Tucker model is a **subspace representation** of a tensor w.r.t. the **tensor basis**¹

$$\{\mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3} : r_i \in [R_i], 1 \leq i \leq 3\},$$

since

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} s_{r_1 r_2 r_3} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3}.$$

Thus, it's natural (and w.l.o.g.) to impose $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$, and similarly for \mathbf{V} and \mathbf{W} .

Indeed, if $\mathbf{X} = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{S}'$, then we can find (semi-)orthogonal \mathbf{U}, \mathbf{V} and \mathbf{W} s.t.

$$(\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{S}' = (\mathbf{U}\mathbf{R}_1, \mathbf{V}\mathbf{R}_2, \mathbf{W}\mathbf{R}_3) \cdot \mathbf{S}' = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot [(\mathbf{R}_1, \mathbf{R}_2, \mathbf{R}_3) \cdot \mathbf{S}'] = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}.$$

In particular, semi-orthogonality of the factors $\mathbf{U}, \mathbf{V}, \mathbf{W}$ implies

$$\mathbf{S} = (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X}.$$

1: Hackbusch, 2012

High-order singular value decomp. (HOSVD)

But even under this constraint, we have a continuum of possible $(\mathbf{U}, \mathbf{V}, \mathbf{W})$, since $\forall (\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3) \in \mathbb{O}(R_1) \times \mathbb{O}(R_2) \times \mathbb{O}(R_3)$,

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} = (\mathbf{U}\mathbf{Q}_1, \mathbf{V}\mathbf{Q}_2, \mathbf{W}\mathbf{Q}_3) \cdot \left[(\mathbf{Q}_1^\top, \mathbf{Q}_2^\top, \mathbf{Q}_3^\top) \cdot \mathbf{S} \right] = (\mathbf{U}', \mathbf{V}', \mathbf{W}') \cdot \mathbf{S}'.$$

A standard choice (useful for low-mrank approximation, as we will see) is:

High-order singular value decomp. (HOSVD)

But even under this constraint, we have a continuum of possible $(\mathbf{U}, \mathbf{V}, \mathbf{W})$, since $\forall (\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3) \in \mathbb{O}(R_1) \times \mathbb{O}(R_2) \times \mathbb{O}(R_3)$,

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} = (\mathbf{U}\mathbf{Q}_1, \mathbf{V}\mathbf{Q}_2, \mathbf{W}\mathbf{Q}_3) \cdot \left[(\mathbf{Q}_1^\top, \mathbf{Q}_2^\top, \mathbf{Q}_3^\top) \cdot \mathbf{S} \right] = (\mathbf{U}', \mathbf{V}', \mathbf{W}') \cdot \mathbf{S}'.$$

A standard choice (useful for low-mrank approximation, as we will see) is:

Def: The **HOSVD** of $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$ is given by¹

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \underset{\text{core tensor}}{\mathbf{S}},$$

where

- $\mathbf{U} \in \mathbb{O}(N_1)$ contains the left singular vectors of $\mathbf{X}_{(1)}$
- $\mathbf{V} \in \mathbb{O}(N_2)$ contains the left singular vectors of $\mathbf{X}_{(2)}$
- $\mathbf{W} \in \mathbb{O}(N_3)$ contains the left singular vectors of $\mathbf{X}_{(3)}$

¹: De Lathauwer & al., 2000a

HOSVD : computation & properties

Computation: Given X , one computes (possibly in parallel) the SVD of each unfolding:

$$\mathbf{X}_{(1)} = \mathbf{U} \Sigma_1 \mathbf{Q}_1^\top,$$

$$\mathbf{X}_{(2)} = \mathbf{V} \Sigma_2 \mathbf{Q}_2^\top,$$

$$\mathbf{X}_{(3)} = \mathbf{W} \Sigma_3 \mathbf{Q}_3^\top,$$

and then $S = (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot X$.

HOSVD : computation & properties

Computation: Given \mathbf{X} , one computes (possibly in parallel) the SVD of each unfolding:

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{\Sigma}_1 \mathbf{Q}_1^\top,$$

$$\mathbf{X}_{(2)} = \mathbf{V} \mathbf{\Sigma}_2 \mathbf{Q}_2^\top,$$

$$\mathbf{X}_{(3)} = \mathbf{W} \mathbf{\Sigma}_3 \mathbf{Q}_3^\top,$$

and then $\mathbf{S} = (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X}$.

Core slice norms: Since $\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top$, we see that

$$\mathbf{S}_{(1)} = \mathbf{\Sigma}_1 \mathbf{Q}_1^\top (\mathbf{W} \boxtimes \mathbf{V}) \quad \Rightarrow \quad \|(\mathbf{S}_{(1)})_{\ell:}\| = (\mathbf{\Sigma}_1)_{\ell\ell} = \|(\mathbf{S})_{\ell::}\|_F$$

and thus $\|(\mathbf{S})_{\ell::}\|_F \geq \|(\mathbf{S})_{\ell+1::}\|_F$. Similarly, $\|(\mathbf{S})_{:\ell}\|_F \geq \|(\mathbf{S})_{:\ell+1:}\|_F$ and $\|(\mathbf{S})_{::\ell}\|_F \geq \|(\mathbf{S})_{::\ell+1}\|_F$.

HOSVD : computation & properties

Computation: Given \mathbf{X} , one computes (possibly in parallel) the SVD of each unfolding:

$$\mathbf{X}_{(1)} = \mathbf{U} \mathbf{\Sigma}_1 \mathbf{Q}_1^\top,$$

$$\mathbf{X}_{(2)} = \mathbf{V} \mathbf{\Sigma}_2 \mathbf{Q}_2^\top,$$

$$\mathbf{X}_{(3)} = \mathbf{W} \mathbf{\Sigma}_3 \mathbf{Q}_3^\top,$$

and then $\mathbf{S} = (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X}$.

Core slice norms: Since $\mathbf{X}_{(1)} = \mathbf{U} \mathbf{S}_{(1)} (\mathbf{W} \boxtimes \mathbf{V})^\top$, we see that

$$\mathbf{S}_{(1)} = \mathbf{\Sigma}_1 \mathbf{Q}_1^\top (\mathbf{W} \boxtimes \mathbf{V}) \quad \Rightarrow \quad \|(\mathbf{S}_{(1)})_{\ell:}\| = (\mathbf{\Sigma}_1)_{\ell\ell} = \|(\mathbf{S})_{\ell::}\|_F$$

and thus $\|(\mathbf{S})_{\ell::}\|_F \geq \|(\mathbf{S})_{\ell+1::}\|_F$. Similarly, $\|(\mathbf{S})_{: \ell}\|_F \geq \|(\mathbf{S})_{: \ell+1}\|_F$ and $\|(\mathbf{S})_{:: \ell}\|_F \geq \|(\mathbf{S})_{:: \ell+1}\|_F$.

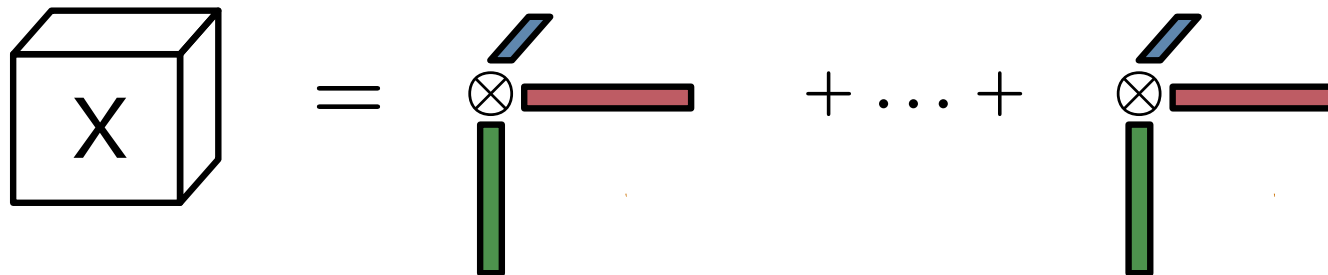
Uniqueness: The HOSVD is unique (up to the signs of the columns of \mathbf{U} , \mathbf{V} , \mathbf{W}) provided that no $\mathbf{X}_{(i)}$ has repeated singular values.

2. The tensor rank

The second notion of rank is defined as:

Def: Rank of X :

$$\text{rank } X := \min \left\{ R \in \mathbb{N} : X = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r, \mathbf{a}_r \in \mathbb{R}^{N_1}, \mathbf{b}_r \in \mathbb{R}^{N_2}, \mathbf{c}_r \in \mathbb{R}^{N_3} \right\}$$



The above set is non-empty: a trivial upper bound is $R \leq N_1 N_2 N_3$, since

$$X = \sum_{n_1, n_2, n_3} (X)_{n_1 n_2 n_3} \underbrace{e_{n_1} \otimes e_{n_2} \otimes e_{n_3}}_{\text{nonzero (= 1) only at } (n_1, n_2, n_3)}.$$

nonzero (= 1) only at (n_1, n_2, n_3)

The tensor rank is a tensor property

Prop: rank X is invariant w.r.t. a change of basis.

Proof: If $X = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$ has rank R , then

$$\text{rank} [(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot X] = \text{rank} \left[\sum_{r=1}^R (\mathbf{U}\mathbf{a}_r) \otimes (\mathbf{V}\mathbf{b}_r) \otimes (\mathbf{W}\mathbf{c}_r) \right] \leq R.$$

It follows that, for $\mathbf{U} \in \text{GL}(N_1)$, $\mathbf{V} \in \text{GL}(N_2)$, $\mathbf{W} \in \text{GL}(N_3)$:

$$\begin{aligned} \text{rank } X &= \text{rank} \{ (\mathbf{U}^{-1}, \mathbf{V}^{-1}, \mathbf{W}^{-1}) \cdot [(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot X] \} \\ &\leq \text{rank}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot X) \\ &\leq \text{rank } X. \end{aligned}$$



More properties will be seen in a moment.

The polyadic decomposition

A polyadic decomposition¹ (PD) of \mathbf{X} of rank R has the form

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R := \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r.$$

1: Hitchcock, 1927

The polyadic decomposition

A polyadic decomposition¹ (PD) of X of rank R has the form

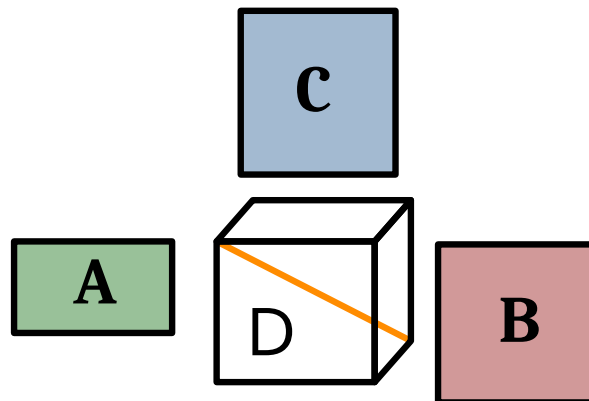
$$X = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R := \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r.$$

It can be seen as a special case of Tucker: by defining

$$\mathbf{D} \in \mathbb{R}^{R \times R \times R}, \quad (\mathbf{D})_{r_1 r_2 r_3} = \delta_{r_1 r_2 r_3},$$

we get

$$(\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D} = \sum_{r_1=1}^R \sum_{r_2=1}^R \sum_{r_3=1}^R (\mathbf{D})_{r_1 r_2 r_3} \mathbf{a}_{r_1} \otimes \mathbf{b}_{r_2} \otimes \mathbf{c}_{r_3} = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$



1: Hitchcock, 1927

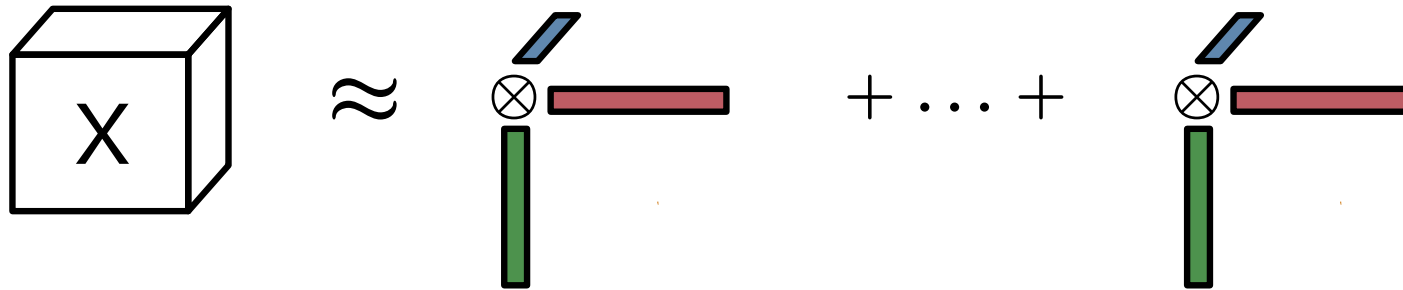
The canonical PD (CPD)

Def: If

$$\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_R = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$

is **minimal** ($\text{rank } \mathbf{X} = R$), then it is called the **CPD** of \mathbf{X} .

It is also often called a **PARAFAC** (for parallel factors) decomposition.



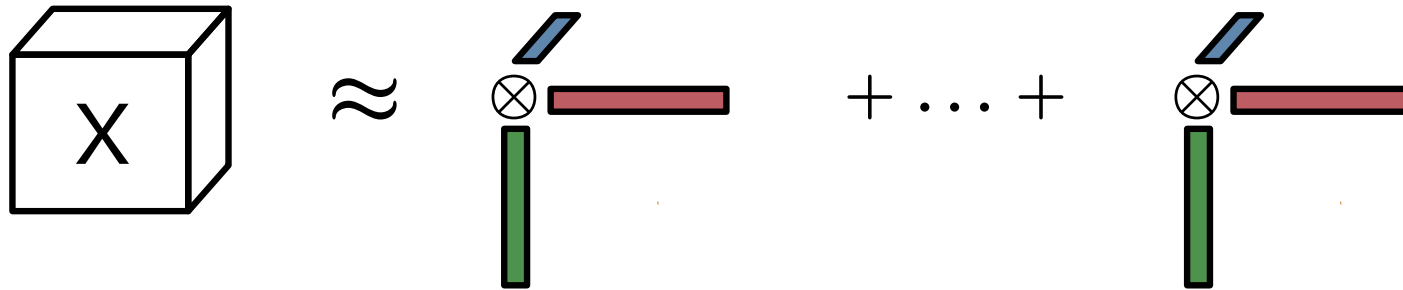
The canonical PD (CPD)

Def: If

$$X = [\mathbf{A}, \mathbf{B}, \mathbf{C}]_R = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$

is **minimal** ($\text{rank } X = R$), then it is called the **CPD** of X .

It is also often called a **PARAFAC** (for parallel factors) decomposition.



Uniqueness: The CPD has quite strong uniqueness properties, discussed ahead.

This key feature has been a major driving force for its study, beginning in the Psychometrics literature.^{1,2,3}

1: Cattell, 1944, 2: Carroll & Chang, 1970, 3: Harshman, 1970

Slices & unfoldings of a CPD

Take $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r \Leftrightarrow x_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}.$

Slices: Frontal slices of \mathbf{X} can be expressed as

$$\mathbb{R}^{N_1 \times N_2} \ni \mathbf{X}_k = (\mathbf{X})_{::k} = \sum_{r=1}^R c_{kr} \mathbf{a}_r \mathbf{b}_r^\top = \mathbf{A} \mathbf{D}_k(\mathbf{C}) \mathbf{B}^\top,$$

with $\mathbf{D}_k(\mathbf{C}) := \text{Diag}(c_{k1}, \dots, c_{kR})$. Similarly, for horizontal and vertical slices:

$$(\mathbf{X})_{i::} = \mathbf{B} \mathbf{D}_k(\mathbf{A}) \mathbf{C}^\top, \quad (\mathbf{X})_{:j:} = \mathbf{A} \mathbf{D}_k(\mathbf{B}) \mathbf{C}^\top.$$

Slices & unfoldings of a CPD

Take $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r \Leftrightarrow x_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}.$

Slices: Frontal slices of \mathbf{X} can be expressed as

$$\mathbb{R}^{N_1 \times N_2} \ni \mathbf{X}_k = (\mathbf{X})_{::k} = \sum_{r=1}^R c_{kr} \mathbf{a}_r \mathbf{b}_r^\top = \mathbf{A} \mathbf{D}_k(\mathbf{C}) \mathbf{B}^\top,$$

with $\mathbf{D}_k(\mathbf{C}) := \text{Diag}(c_{k1}, \dots, c_{kR})$. Similarly, for horizontal and vertical slices:

$$(\mathbf{X})_{i::} = \mathbf{B} \mathbf{D}_k(\mathbf{A}) \mathbf{C}^\top, \quad (\mathbf{X})_{:j:} = \mathbf{A} \mathbf{D}_k(\mathbf{B}) \mathbf{C}^\top.$$

Unfoldings: By linearity of unfolding,

$$\mathbf{X}_{(1)} = \sum_{r=1}^R (\mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r)_{(1)} = \sum_{r=1}^R \mathbf{a}_r (\mathbf{c}_r \boxtimes \mathbf{b}_r)^\top = \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top,$$

Def: [Khatri-Rao](#) product (columnwise \boxtimes): $\mathbf{C} \odot \mathbf{B} := (\mathbf{c}_1 \boxtimes \mathbf{b}_1 \quad \dots \quad \mathbf{c}_R \boxtimes \mathbf{b}_R).$

Unfoldings & vectorization of a CPD

Unfoldings: By the same reasoning,

$$\mathbf{X}_{(2)} = \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top, \quad \mathbf{X}_{(3)} = \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top.$$

Unfoldings & vectorization of a CPD

Unfoldings: By the same reasoning,

$$\mathbf{X}_{(2)} = \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top, \quad \mathbf{X}_{(3)} = \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top.$$

Vectorization: Again by linearity,

$$\text{vec } \mathbf{X} = \sum_{r=1}^R \text{vec}(\mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r) = \sum_{r=1}^R \mathbf{c}_r \boxtimes \mathbf{b}_r \boxtimes \mathbf{a}_r = (\mathbf{C} \odot \mathbf{B} \odot \mathbf{A})\mathbf{1}.$$

Computation of a CPD

In special cases, a CPD can be computed by non-iterative algebraic methods which rely on standard numerical linear algebra routines.

However, no general algorithm of that kind exists.

Computation of a CPD

In special cases, a CPD can be computed by non-iterative **algebraic methods** which rely on standard numerical linear algebra routines.

However, no general algorithm of that kind exists.

Typically, the problem is instead addressed by fitting the CPD model to the data according to some criterion or loss function.

Most usual approach: given a rank R , solve

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2.$$

Nonconvex problem, since the model is multilinear in \mathbf{A} , \mathbf{B} , \mathbf{C} .

Computation of a CPD

In special cases, a CPD can be computed by non-iterative **algebraic methods** which rely on standard numerical linear algebra routines.

However, no general algorithm of that kind exists.

Typically, the problem is instead addressed by fitting the CPD model to the data according to some criterion or loss function.

Most usual approach: given a rank R , solve

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2.$$

Nonconvex problem, since the model is multilinear in \mathbf{A} , \mathbf{B} , \mathbf{C} .

Clearly, this seems well-suited for approximation as well, which is the most common objective in applications.

However, low-rank approximation comes with its own pitfalls, as we'll see.

Alternating least squares (ALS)

Assume $\mathbf{X} = \llbracket \mathbf{A}^*, \mathbf{B}^*, \mathbf{C}^* \rrbracket_R$. How can we address $\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2$?

Alternating least squares (ALS)

Assume $\mathbf{X} = \llbracket \mathbf{A}^*, \mathbf{B}^*, \mathbf{C}^* \rrbracket_R$. How can we address $\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2$?

The simplest and best known algorithm is a BCD scheme called **alternating least squares (ALS)**.^{1,2}

Idea: While the original problem is (jointly) nonconvex, it is block-convex:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2 = \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \right\|_F^2 \quad (\text{i})$$

$$= \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top \right\|_F^2 \quad (\text{ii})$$

$$= \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top \right\|_F^2 \quad (\text{iii})$$

1: Harshman, 1970, 2: Carroll & Chang, 1970

Alternating least squares (ALS)

Assume $\mathbf{X} = \llbracket \mathbf{A}^*, \mathbf{B}^*, \mathbf{C}^* \rrbracket_R$. How can we address $\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2$?

The simplest and best known algorithm is a BCD scheme called **alternating least squares (ALS)**.^{1,2}

Idea: While the original problem is (jointly) nonconvex, it is block-convex:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2 = \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \right\|_F^2 \quad (\text{i})$$

$$= \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(2)} - \mathbf{B}(\mathbf{C} \odot \mathbf{A})^\top \right\|_F^2 \quad (\text{ii})$$

$$= \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X}_{(3)} - \mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top \right\|_F^2 \quad (\text{iii})$$

Hence, conditioned upon the current estimates \mathbf{B} and \mathbf{C} , the solution for \mathbf{A} can be explicitly derived from (i) as (assuming $\text{rank } \mathbf{C} \odot \mathbf{B} = R$):

$$\begin{aligned} \mathbf{A} &\leftarrow \mathbf{X}_{(1)} [(\mathbf{C} \odot \mathbf{B})^\dagger]^\top = \mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) [(\mathbf{C} \odot \mathbf{B})^\top (\mathbf{C} \odot \mathbf{B})]^{-1} \\ &= \mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) [(\mathbf{C}^\top \mathbf{C}) \circ (\mathbf{B}^\top \mathbf{B})]^{-1} \end{aligned}$$

1: Harshman, 1970, 2: Carroll & Chang, 1970

ALS: properties and variants

In its “vanilla” form, ALS is prone to facing ill-conditioning issues, slow convergence or failure to converge in reasonable time.

Nevertheless, ALS & its variants are wildly popular, as they are simple and often exhibit a good performance.

ALS: properties and variants

In its “vanilla” form, ALS is prone to facing ill-conditioning issues, slow convergence or failure to converge in reasonable time.

Nevertheless, ALS & its variants are wildly popular, as they are simple and often exhibit a good performance.

Some common variants are:

- “enhanced line search” (ELS) methods for accelerating convergence¹
- adding a proximal term (w.r.t. current estimate \mathbf{A}_0) to each subproblem to circumvent ill-conditioning & guarantee well-posed updates:

$$\min_{\mathbf{A}} \left\| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T \right\|_F^2 + \tau \|\mathbf{A} - \mathbf{A}_0\|_F^2$$

(strictly convex problem, with $2\tau\mathbf{I}$ added to the Hessian)

1: Rajih & al., 2008

ALS: convergence

By construction, the objective value is nonincreasing along ALS iterations.
But what about convergence of iterates?

ALS: convergence

By construction, the objective value is nonincreasing along ALS iterations.

But what about convergence of iterates?

Local convergence: (Uschmajew, 2012) showed

- local convergence under non-degeneracy condition of local min
- local linear convergence around global minimizers of

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]_R\|_F^2 + \tau \left(\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 + \|\mathbf{C}\|_F^2 \right)$$

ALS: convergence

By construction, the objective value is nonincreasing along ALS iterations.
But what about convergence of iterates?

Local convergence: (Uschmajew, 2012) showed

- local convergence under non-degeneracy condition of local min
- local linear convergence around global minimizers of

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R\|_F^2 + \tau \left(\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 + \|\mathbf{C}\|_F^2 \right)$$

Global convergence:

- Implied by the general framework of (Xu & al., 2013) when

$$\lambda_{\min} \left((\mathbf{C} \odot \mathbf{B})^\top (\mathbf{C} \odot \mathbf{B}) \right) \geq c, \quad \text{for some } c > 0$$

at every iteration, and idem for $(\mathbf{C} \odot \mathbf{A})^\top (\mathbf{C} \odot \mathbf{A})$ and $(\mathbf{B} \odot \mathbf{A})^\top (\mathbf{B} \odot \mathbf{A})$

- (Yang, 2023) proposed a variant which converges if we instead bound the **smallest positive eigenvalues** λ_{\min}^+ of those matrices uniformly

Algebraic sol'n by simultaneous diagonalization

A common method for the initialization of ALS (or any other iterative algorithm) is by simultaneous diagonalization (or “Jennrich’s algorithm,” a misnomer).^{1,2}

It can be used for $d = 3$ whenever, say, $\text{rank } \mathbf{X} = R \leq \min \{N_1, N_2\}$ w.l.o.g.

As we’ll see, $\text{mrnk } \mathbf{X} \leq R$ entry-wise, thus one can project \mathbf{X} onto $\mathbb{R}^{R \times R \times N_3}$. Hence, w.l.o.g. we assume $\mathbf{X} \in \mathbb{R}^{R \times R \times N_3}$.

1: Sanchez & Kowalski, 1990, 2: Leurgans & al., 1993

Algebraic sol'n by simultaneous diagonalization

A common method for the initialization of ALS (or any other iterative algorithm) is by simultaneous diagonalization (or “Jennrich’s algorithm,” a misnomer).^{1,2}

It can be used for $d = 3$ whenever, say, $\text{rank } \mathbf{X} = R \leq \min \{N_1, N_2\}$ w.l.o.g.

As we’ll see, $\text{mrnk } \mathbf{X} \leq R$ entry-wise, thus one can project \mathbf{X} onto $\mathbb{R}^{R \times R \times N_3}$. Hence, w.l.o.g. we assume $\mathbf{X} \in \mathbb{R}^{R \times R \times N_3}$.

Hypothesis: $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R$, with $\text{rank } \mathbf{A} = \text{rank } \mathbf{B} = R$ and \mathbf{C} doesn’t have collinear columns.

Idea: By computing contractions $\mathbf{Y}_1 = (\cdot, \cdot, \boldsymbol{\theta}_1) \cdot \mathbf{X}$ and $\mathbf{Y}_2 = (\cdot, \cdot, \boldsymbol{\theta}_2) \cdot \mathbf{X}$, we get

$$\mathbf{Y}_i = \sum_{r=1}^R (c_r^\top \boldsymbol{\theta}_i) a_r b_r^\top = \mathbf{A} \text{Diag}(c_1^\top \boldsymbol{\theta}_i, \dots, c_R^\top \boldsymbol{\theta}_i) \mathbf{B}^\top = \mathbf{A} \mathbf{D}_i \mathbf{B}^\top$$

Suppose \mathbf{D}_2 is nonsingular (generically true for random $\boldsymbol{\theta}_2$). Then,

$$\mathbf{Y} := \mathbf{Y}_1 \mathbf{Y}_2^{-1} = \mathbf{A} \mathbf{D}_1 \mathbf{D}_2^{-1} \mathbf{A}^{-1}$$

so that the eigenvectors of \mathbf{Y} are cols of \mathbf{A} (with distinct eigenvalues)!

1: Sanchez & Kowalski, 1990, 2: Leurgans & al., 1993

Drawbacks, properties & extensions

Once \mathbf{A} is recovered, one can compute \mathbf{B} and \mathbf{C} from

$$(\mathbf{A}^{-1}\mathbf{X}_{(1)})^{\top} = \mathbf{C} \odot \mathbf{B},$$

by rank-1 approximation of each column.

Drawbacks, properties & extensions

Once \mathbf{A} is recovered, one can compute \mathbf{B} and \mathbf{C} from

$$(\mathbf{A}^{-1}\mathbf{X}_{(1)})^{\top} = \mathbf{C} \odot \mathbf{B},$$

by rank-1 approximation of each column.

Drawbacks: Possible numerical instability¹ and lack of robustness to noise.

1: Beltrán & al., 2019, 2: Bhaskara & al., 2014

Drawbacks, properties & extensions

Once \mathbf{A} is recovered, one can compute \mathbf{B} and \mathbf{C} from

$$(\mathbf{A}^{-1} \mathbf{X}_{(1)})^{\top} = \mathbf{C} \odot \mathbf{B},$$

by rank-1 approximation of each column.

Drawbacks: Possible numerical instability¹ and lack of robustness to noise.

Analysis: Error bounds were given by Bhaskara & al. (2014) under a “smoothed model,” under additive noise with bounded entries (depending on the conditioning of \mathbf{A} , \mathbf{B} and the coherence of \mathbf{C}).

1: Beltrán & al., 2019, 2: Bhaskara & al., 2014

Drawbacks, properties & extensions

Once \mathbf{A} is recovered, one can compute \mathbf{B} and \mathbf{C} from

$$(\mathbf{A}^{-1} \mathbf{X}_{(1)})^\top = \mathbf{C} \odot \mathbf{B},$$

by rank-1 approximation of each column.

Drawbacks: Possible numerical instability¹ and lack of robustness to noise.

Analysis: Error bounds were given by Bhaskara & al. (2014) under a “smoothed model,” under additive noise with bounded entries (depending on the conditioning of \mathbf{A} , \mathbf{B} and the coherence of \mathbf{C}).

Extensions: Higher-order tensors can be decomposed by partially ignoring the tensor structure. For instance, if $d = 2m + 1$ with $m > 1$, then one takes

$$\sum_{r=1}^R \mathbf{a}^{(1)} \otimes \dots \otimes \mathbf{a}^{(d)} \mapsto \sum_{r=1}^R (\mathbf{a}^{(1)} \boxtimes \dots \boxtimes \mathbf{a}^{(m)}) \otimes (\mathbf{a}^{(m+1)} \boxtimes \dots \boxtimes \mathbf{a}^{(2m)}) \otimes \mathbf{a}^{(d)}.$$

Note that now R can exceed the dimensions N_i of all factors.

1: Beltrán & al., 2019, 2: Bhaskara & al., 2014

Properties of the tensor rank, part II

(i) Lower bound by mrank:

Prop: rank X is lower bounded by each component of mrank X .

Proof: Let rank $X = R$ and write the CPD $X = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$. Then,

$$\text{rank}_1 X = \text{rank } \mathbf{X}_{(1)} = \text{rank } \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \leq R,$$

and similarly for the other modes. ■

Properties of the tensor rank, part II

(i) Lower bound by mrank:

Prop: rank X is lower bounded by each component of mrank X .

Proof: Let rank $X = R$ and write the CPD $X = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$. Then,

$$\text{rank}_1 X = \text{rank } \mathbf{X}_{(1)} = \text{rank } \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \leq R,$$

and similarly for the other modes. ■

Rmk: This implies that rank $X \geq \text{rank } \mathbf{X}_k$ for all k (similarly for horizontal and vertical slices).

Properties of the tensor rank, part II

(i) Lower bound by mrank:

Prop: rank X is lower bounded by each component of mrank X .

Proof: Let rank $X = R$ and write the CPD $X = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$. Then,

$$\text{rank}_1 X = \text{rank } \mathbf{X}_{(1)} = \text{rank } \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top \leq R,$$

and similarly for the other modes. ■

Rmk: This implies that rank $X \geq \text{rank } \mathbf{X}_k$ for all k (similarly for horizontal and vertical slices).

(ii) Rank invariance under embedding in larger space:

Prop: If $X = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{S}$ and $\mathbf{A}, \mathbf{B}, \mathbf{C}$ have full column rank, then rank $X = \text{rank } \mathbf{S}$.

Proof: See exercices.

Rmk: Often used (in an approximate fashion) when computing the CPD of a low-mrank tensor: one can work in $\mathbb{R}^{R_1 \times R_2 \times R_3}$ instead of in $\mathbb{R}^{N_1 \times N_2 \times N_3}$

Tensor rank: examples

Example: If

$$\mathbf{D} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \cdots + \mathbf{e}_R \otimes \mathbf{e}_R \otimes \mathbf{e}_R,$$

then $\text{rank } \mathbf{D} = R$, since $\text{rank } \mathbf{D} \geq \text{rank}_1 \mathbf{D} = \text{rank } \mathbf{I} (\mathbf{I} \odot \mathbf{I})^\top = R$.



Tensor rank: examples

Example: If

$$\mathbf{D} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \cdots + \mathbf{e}_R \otimes \mathbf{e}_R \otimes \mathbf{e}_R,$$

then $\text{rank } \mathbf{D} = R$, since $\text{rank } \mathbf{D} \geq \text{rank}_1 \mathbf{D} = \text{rank } \mathbf{I} (\mathbf{I} \odot \mathbf{I})^\top = R$. □

Example: Let \mathbf{A} , \mathbf{B} , \mathbf{C} have full column rank, and define

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r.$$

We have

$$\mathbf{X} = \sum_{r=1}^R (\mathbf{A} \mathbf{e}_r) \otimes (\mathbf{B} \mathbf{e}_r) \otimes (\mathbf{C} \mathbf{e}_r) = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D},$$

thus $\text{rank } \mathbf{X} = R$. □

Tensor rank: examples

Example: If

$$\mathbf{D} = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \cdots + \mathbf{e}_R \otimes \mathbf{e}_R \otimes \mathbf{e}_R,$$

then $\text{rank } \mathbf{D} = R$, since $\text{rank } \mathbf{D} \geq \text{rank}_1 \mathbf{D} = \text{rank } \mathbf{I} (\mathbf{I} \odot \mathbf{I})^\top = R$. □

Example: Let $\mathbf{A}, \mathbf{B}, \mathbf{C}$ have full column rank, and define

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r.$$

We have

$$\mathbf{X} = \sum_{r=1}^R (\mathbf{A} \mathbf{e}_r) \otimes (\mathbf{B} \mathbf{e}_r) \otimes (\mathbf{C} \mathbf{e}_r) = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D},$$

thus $\text{rank } \mathbf{X} = R$. □

Example:

$$\text{rank} (\mathbf{a}_1 \otimes \mathbf{b} \otimes \mathbf{c} + \mathbf{a}_2 \otimes \mathbf{b} \otimes \mathbf{c}) = \text{rank} (\mathbf{a}_1 + \mathbf{a}_2) \otimes \mathbf{b} \otimes \mathbf{c} \leq 1,$$

while

$$\text{rank} (\mathbf{a} \otimes \mathbf{b}_1 \otimes \mathbf{c}_1 + \mathbf{a} \otimes \mathbf{b}_2 \otimes \mathbf{c}_2) \leq 2. \quad \square$$

Properties of the tensor rank, part III

(iii) Upper bound by product of mrnk components:

Prop: If $\text{mrnk } \mathbf{X} = (R_1, R_2, R_3)$, then $\text{rank } \mathbf{X} \leq \min \{R_1 R_2, R_1 R_3, R_2 R_3\}$.

Proof: From an $\text{mrnk}-(R_1, R_2, R_3)$ Tucker representation of \mathbf{X} , we get

$$\mathbf{X} = \sum_{r_1, r_2, r_3} s_{r_1 r_2 r_3} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3} = \sum_{r_1, r_2} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \left(\sum_{r_3} s_{r_1 r_2 r_3} \mathbf{w}_{r_3} \right),$$

which shows that $\text{rank } \mathbf{X} \leq R_1 R_2$. The other bounds are shown similarly. ■

Properties of the tensor rank, part III

(iii) Upper bound by product of mrnk components:

Prop: If $\text{mrnk } \mathbf{X} = (R_1, R_2, R_3)$, then $\text{rank } \mathbf{X} \leq \min \{R_1 R_2, R_1 R_3, R_2 R_3\}$.

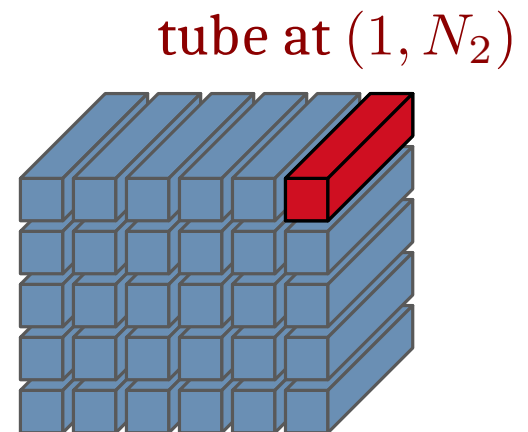
Proof: From an $\text{mrnk}-(R_1, R_2, R_3)$ Tucker representation of \mathbf{X} , we get

$$\mathbf{X} = \sum_{r_1, r_2, r_3} s_{r_1 r_2 r_3} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \mathbf{w}_{r_3} = \sum_{r_1, r_2} \mathbf{u}_{r_1} \otimes \mathbf{v}_{r_2} \otimes \left(\sum_{r_3} s_{r_1 r_2 r_3} \mathbf{w}_{r_3} \right),$$

which shows that $\text{rank } \mathbf{X} \leq R_1 R_2$. The other bounds are shown similarly. ■

A more intuitive way of seeing this: assume w.l.o.g. that $\text{mrnk } \mathbf{X} = (N_1, N_2, N_3)$ and “synthetize” \mathbf{X} as

$$\begin{aligned} \mathbf{X} &= \sum_{n_1, n_2, n_3} x_{n_1 n_2 n_3} \mathbf{e}_{n_1} \otimes \mathbf{e}_{n_2} \otimes \mathbf{e}_{n_3} \\ &= \sum_{n_1, n_2} \mathbf{e}_{n_1} \otimes \mathbf{e}_{n_2} \otimes \underbrace{\left(\sum_{n_3} x_{n_1 n_2 n_3} \mathbf{e}_{n_3} \right)}_{\text{“tube” at } (n_1, n_2)}. \end{aligned}$$



Properties of the tensor rank, part III cont'd

(iv) rank X depends on the underlying field

Example: Take

$$\begin{aligned}x_{ijk} &= \cos(i + j + k) = \frac{1}{2} \left(e^{i(i+j+k)} + e^{-i(i+j+k)} \right) \\&= \frac{1}{2} e^{ii} e^{ij} e^{ik} + \frac{1}{2} \overline{e^{ii} e^{ij} e^{ik}}.\end{aligned}$$

Hence, rank $X = 2$ over \mathbb{C} .

Properties of the tensor rank, part III cont'd

(iv) rank X depends on the underlying field

Example: Take

$$\begin{aligned} x_{ijk} &= \cos(i + j + k) = \frac{1}{2} \left(e^{i(i+j+k)} + e^{-i(i+j+k)} \right) \\ &= \frac{1}{2} e^{ii} e^{ij} e^{ik} + \frac{1}{2} \overline{e^{ii} e^{ij} e^{ik}}. \end{aligned}$$

Hence, rank $X = 2$ over \mathbb{C} .

By using again Euler's identity $e^{i\ell} = a_\ell + ib_\ell$ with $a_\ell := \cos \ell$ and $b_\ell := \sin \ell$,

$$\begin{aligned} x_{ijk} &= a_i a_j a_k - a_i b_j b_k - b_i a_j b_k - b_i b_j a_k \quad (0 \text{ or } 2 \text{ factors } b) \\ &= \left[(\mathbf{H}, \mathbf{H}, \mathbf{H}) \cdot \underbrace{\begin{pmatrix} 1 & 0 & | & 0 & -1 \\ 0 & -1 & | & -1 & 0 \end{pmatrix}}_{\text{rank}=3} \right]_{ijk}, \quad \text{with } \mathbf{H} := \begin{pmatrix} a & b \end{pmatrix} \\ &= -(a_i + b_i) b_j b_k + (a_i - b_i) a_j a_k + b_i (a_j - b_j) (a_k - b_k). \end{aligned}$$

Hence, rank $X = 3$ over \mathbb{R} !



Properties of the tensor rank, part III cont'd

(v) rank X can exceed **all** tensor dimensions

Example: Let

$$X = \left(\begin{array}{cc|cc} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{array} \right) = e_1 \otimes e_1 \otimes e_1 + e_2 \otimes e_2 \otimes e_1 + e_1 \otimes e_2 \otimes e_2.$$

Is this minimal?

Properties of the tensor rank, part III cont'd

(v) rank X can exceed **all** tensor dimensions

Example: Let

$$X = \left(\begin{array}{cc|cc} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{array} \right) = e_1 \otimes e_1 \otimes e_1 + e_2 \otimes e_2 \otimes e_1 + e_1 \otimes e_2 \otimes e_2.$$

Is this minimal?

First, note that $\text{rank } X \geq \text{rank}_1 X = 2$. Suppose that $\text{rank } X = 2$. Then, we can write $X = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ with $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{2 \times 2}$ and

$$\mathbf{X}_1 = \mathbf{I} = \mathbf{A}\mathbf{D}_1(\mathbf{C})\mathbf{B}^\top, \quad \mathbf{X}_2 = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = \mathbf{A}\mathbf{D}_2(\mathbf{C})\mathbf{B}^\top.$$

Properties of the tensor rank, part III cont'd

(v) rank X can exceed **all** tensor dimensions

Example: Let

$$X = \left(\begin{array}{cc|cc} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{array} \right) = e_1 \otimes e_1 \otimes e_1 + e_2 \otimes e_2 \otimes e_1 + e_1 \otimes e_2 \otimes e_2.$$

Is this minimal?

First, note that $\text{rank } X \geq \text{rank}_1 X = 2$. Suppose that $\text{rank } X = 2$. Then, we can write $X = \llbracket A, B, C \rrbracket$ with $A, B, C \in \mathbb{R}^{2 \times 2}$ and

$$X_1 = I = AD_1(C)B^T, \quad X_2 = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = AD_2(C)B^T.$$

As X_1 is nonsingular, it follows that

$$X_2 X_1^{-1} = X_2 = X_2 B^{-T} (D_1(C))^{-1} A^{-1} = A \underbrace{D_2(C) (D_1(C))^{-1}}_{:=D} A^{-1},$$

implying X_2 is diagonalizable—a contradiction!

Conclusion: **rank** $X = 3$. This is actually the maximal rank on $\mathbb{R}^{2 \times 2 \times 2}$.



Maximal rank

Q: What is the maximal rank of a tensor from $\mathbb{R}^{N_1 \times N_2 \times N_3}$?

Many results have been derived for more or less particular cases, e.g.:

- maximal rank over $\mathbb{R}^{2 \times 2 \times 2}$ or $\mathbb{C}^{2 \times 2 \times 2}$ is¹ 3
- maximal rank over $\mathbb{R}^{3 \times 3 \times 3}$ or $\mathbb{C}^{3 \times 3 \times 3}$ is^{1,2} 5
- maximal rank over $\mathbb{R}^{N_1 \times N_2 \times N_3}$ is³ at most $N_1 + \lfloor N_3/2 \rfloor N_2$

See the references given in the handout for more results.

This question is not of much concern for the applications that we consider here: both uniqueness and compressibility require a much smaller rank.

1: Kruskal, 1989, 2: Bremner & Hu, 2013, 3: Sumi & al., 2010

Generic and typical ranks

Recap for matrices:

Let $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$ be such that $x_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. **Q:** What is rank \mathbf{X} ?

Generic and typical ranks

Recap for matrices:

Let $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$ be such that $x_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. **Q:** What is rank \mathbf{X} ?

Answer: $\mathbb{P}(\text{rank } \mathbf{X} = \min \{N_1, N_2\}) = 1$, $\mathbb{P}(\text{rank } \mathbf{X} < \min \{N_1, N_2\}) = 0$

The set of rank-defective matrices has null Lebesgue measure \Rightarrow any absolutely continuous distribution (w.r.t. Lebesgue) will assign zero mass to it.

What about tensors?

Generic and typical ranks

Recap for matrices:

Let $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2}$ be such that $x_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. **Q:** What is rank \mathbf{X} ?

Answer: $\mathbb{P}(\text{rank } \mathbf{X} = \min \{N_1, N_2\}) = 1$, $\mathbb{P}(\text{rank } \mathbf{X} < \min \{N_1, N_2\}) = 0$

The set of rank-defective matrices has null Lebesgue measure \Rightarrow any absolutely continuous distribution (w.r.t. Lebesgue) will assign zero mass to it.

What about tensors?

Example¹: Let's do a similar experiment: $\mathbf{X} \in \mathbb{R}^{2 \times 2 \times 2}$, with $x_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. Then²:

- over \mathbb{R} , $\mathbb{P}(\text{rank } \mathbf{X} = 2) = \frac{\pi}{4}$, $\mathbb{P}(\text{rank } \mathbf{X} = 3) = 1 - \frac{\pi}{4}$
- over \mathbb{C} , $\mathbb{P}(\text{rank } \mathbf{X} = 2) = 1$

What is going on?

1: Sidiropoulos & al., 2017, 2: Bergqvist, 2013

Generic and typical ranks, cont'd

Write $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ the CPD of \mathbf{X} . Since $\text{rank } \mathbf{X}_1 = \text{rank } \mathbf{X}_2 = 2$ a.s., if $\text{rank } \mathbf{X} = 2$ over \mathbb{R} , then all (real-valued) matrices in

$$\mathbf{X}_1 = \mathbf{A} \mathbf{D}_1(\mathbf{C}) \mathbf{B}^\top, \quad \mathbf{X}_2 = \mathbf{A} \mathbf{D}_2(\mathbf{C}) \mathbf{B}^\top$$

must be nonsingular.

Generic and typical ranks, cont'd

Write $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ the CPD of \mathbf{X} . Since $\text{rank } \mathbf{X}_1 = \text{rank } \mathbf{X}_2 = 2$ a.s., if $\text{rank } \mathbf{X} = 2$ over \mathbb{R} , then all (real-valued) matrices in

$$\mathbf{X}_1 = \mathbf{A} \mathbf{D}_1(\mathbf{C}) \mathbf{B}^\top, \quad \mathbf{X}_2 = \mathbf{A} \mathbf{D}_2(\mathbf{C}) \mathbf{B}^\top$$

must be nonsingular. Next, we compute (does this look familiar?):

$$\mathbf{X}_1 \mathbf{X}_2^{-1} = \mathbf{A} \mathbf{D}_1(\mathbf{C}) (\mathbf{D}_2(\mathbf{C}))^{-1} \mathbf{A}^{-1} = \mathbf{A} \mathbf{D} \mathbf{A}^{-1}.$$

But real-valued matrices can have complex-valued eigenvalues with positive probability, in which case we would have a contradiction.

By contrast, this works with probability one if $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{C}^{2 \times 2}$, hence $\text{rank } \mathbf{X} = 2$ a.s. over \mathbb{C} . □

Generic and typical ranks, cont'd

Write $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ the CPD of \mathbf{X} . Since $\text{rank } \mathbf{X}_1 = \text{rank } \mathbf{X}_2 = 2$ a.s., if $\text{rank } \mathbf{X} = 2$ over \mathbb{R} , then all (real-valued) matrices in

$$\mathbf{X}_1 = \mathbf{A} \mathbf{D}_1(\mathbf{C}) \mathbf{B}^\top, \quad \mathbf{X}_2 = \mathbf{A} \mathbf{D}_2(\mathbf{C}) \mathbf{B}^\top$$

must be nonsingular. Next, we compute (does this look familiar?):

$$\mathbf{X}_1 \mathbf{X}_2^{-1} = \mathbf{A} \mathbf{D}_1(\mathbf{C}) (\mathbf{D}_2(\mathbf{C}))^{-1} \mathbf{A}^{-1} = \mathbf{A} \mathbf{D} \mathbf{A}^{-1}.$$

But real-valued matrices can have complex-valued eigenvalues with positive probability, in which case we would have a contradiction.

By contrast, this works with probability one if $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{C}^{2 \times 2}$, hence $\text{rank } \mathbf{X} = 2$ a.s. over \mathbb{C} . □

In general: For tensors in $\mathbb{R}^{N_1 \times \cdots \times N_d}$,

- there is a unique¹ **generic rank** over \mathbb{C} , smaller than the maximal rank
- there are multiple **typical ranks** in \mathbb{R} (the smallest one is generic¹ over \mathbb{C}).

How large is the generic rank?

For an $N_1 \times \cdots \times N_d$ tensor over \mathbb{C} , the generic rank is, according to the Abo–Ottaviani–Peterson conjecture^{1,2}

$$\text{grank}(N_1, \dots, N_d) = \left\lceil \frac{\prod_{i=1}^d N_i}{1 + \sum_{i=1}^d (N_i - 1)} \right\rceil,$$

with some exceptions.

As we will discuss ahead, CPD uniqueness generically holds up to $\text{grank} - 1$ (but not beyond).

1: Abo & al., 2009, 2: Vannieuwenhoven, 2015

The block-term decomposition (BTD) model

What happens if two or more terms of a PD share a vector? For example:

$$\begin{aligned} X &= \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c} + \mathbf{a}_2 \otimes \mathbf{b}_2 \otimes \mathbf{c} + \dots \\ &= (\mathbf{a}_1 \otimes \mathbf{b}_1 + \mathbf{a}_2 \otimes \mathbf{b}_2) \otimes \mathbf{c} + \dots \\ &= (\mathbf{AB}^\top) \otimes \mathbf{c} + \dots \end{aligned}$$

“Full” identifiability breaks down, since \mathbf{A} and \mathbf{B} are not identifiable from \mathbf{AB}^\top .

But can at least \mathbf{AB}^\top and \mathbf{c} be identifiable?

The block-term decomposition (BTD) model

What happens if two or more terms of a PD share a vector? For example:

$$\begin{aligned} X &= \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c} + \mathbf{a}_2 \otimes \mathbf{b}_2 \otimes \mathbf{c} + \dots \\ &= (\mathbf{a}_1 \otimes \mathbf{b}_1 + \mathbf{a}_2 \otimes \mathbf{b}_2) \otimes \mathbf{c} + \dots \\ &= (\mathbf{AB}^\top) \otimes \mathbf{c} + \dots \end{aligned}$$

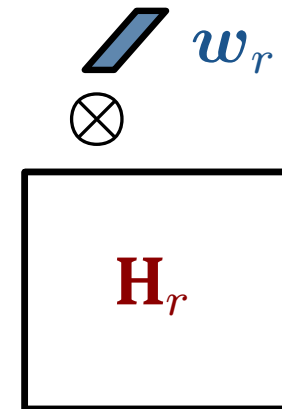
“Full” identifiability breaks down, since \mathbf{A} and \mathbf{B} are not identifiable from \mathbf{AB}^\top .

But can at least \mathbf{AB}^\top and \mathbf{c} be identifiable?

Yes, under a low-rank constraint.

This is the point of the BTD model:¹

$$X = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{c}_r, \quad \text{rank } \mathbf{H}_r \leq L_r$$



Various uniqueness results also exist, some of which we will discuss ahead.

1: De Lathauwer, 2008

BTD as a sum of low-mrank tensors

Q: If $\mathbf{X} = \mathbf{H} \otimes \mathbf{c}$ with $\text{rank } \mathbf{H} = L$, then $\text{mrnk } \mathbf{X} = ?$

Write

$$\mathbf{X}_{(1)} = (c_1 \mathbf{H} \quad c_2 \mathbf{H} \quad \dots \quad c_{N_3} \mathbf{H}) = (\mathbf{1}_{N_3}^\top \boxtimes \mathbf{H}) \text{Diag}(\mathbf{c} \boxtimes \mathbf{1}_{N_2})$$

to see that $\text{rank}_1 \mathbf{X} = L$ (similarly for rank_2). Also,

$$\mathbf{X}_{(3)} = \mathbf{c} \text{vec}(\mathbf{H})^\top \Rightarrow \text{rank}_3 \mathbf{X} = 1.$$

Hence, $\text{mrnk } \mathbf{X} = (L, L, 1)$ & we can see the BTD as a sum of low-mrank blocks.

BTD as a sum of low-mrank tensors

Q: If $\mathbf{X} = \mathbf{H} \otimes \mathbf{c}$ with $\text{rank } \mathbf{H} = L$, then $\text{mrnk } \mathbf{X} = ?$

Write

$$\mathbf{X}_{(1)} = (c_1 \mathbf{H} \quad c_2 \mathbf{H} \quad \dots \quad c_{N_3} \mathbf{H}) = (\mathbf{1}_{N_3}^\top \boxtimes \mathbf{H}) \text{Diag}(\mathbf{c} \boxtimes \mathbf{1}_{N_2})$$

to see that $\text{rank}_1 \mathbf{X} = L$ (similarly for rank_2). Also,

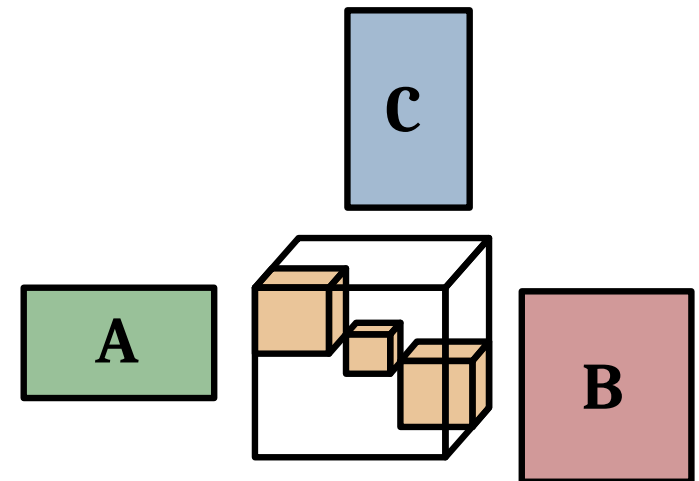
$$\mathbf{X}_{(3)} = \mathbf{c} \text{vec}(\mathbf{H})^\top \Rightarrow \text{rank}_3 \mathbf{X} = 1.$$

Hence, $\text{mrnk } \mathbf{X} = (L, L, 1)$ & we can see the BTD as a sum of low-mrank blocks.

More generally, it can be defined that way:¹

$$\mathbf{X} = \sum_{r=1}^R (\mathbf{A}_r, \mathbf{B}_r, \mathbf{C}_r) \cdot \mathbf{S}_r$$

Some uniqueness results also exist for this more general model.¹



¹: De Lathauwer, 2008

BTD computation

Just like the CPD, if \mathbf{X} is given by a BTD

$$\mathbf{X} = \sum_{r=1}^R [\mathbf{A}_r^* (\mathbf{B}_r^*)^\top] \otimes \mathbf{c}_r^*,$$

with $\mathbf{A}_r^* \in \mathbb{R}^{N_1 \times L_r}$, $\mathbf{B}_r^* \in \mathbb{R}^{N_2 \times L_r}$ and $\mathbf{c}_r^* \in \mathbb{R}^{N_3}$, then its computation can be carried out via

$$\min_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{c}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{c}_r \right\|_{\text{F}}^2,$$

with same dimensions.

BTD computation

Just like the CPD, if \mathbf{X} is given by a BTD

$$\mathbf{X} = \sum_{r=1}^R [\mathbf{A}_r^* (\mathbf{B}_r^*)^\top] \otimes \mathbf{c}_r^*,$$

with $\mathbf{A}_r^* \in \mathbb{R}^{N_1 \times L_r}$, $\mathbf{B}_r^* \in \mathbb{R}^{N_2 \times L_r}$ and $\mathbf{c}_r^* \in \mathbb{R}^{N_3}$, then its computation can be carried out via

$$\min_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{c}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{c}_r \right\|_F^2,$$

with same dimensions.

In particular, an ALS algorithm can also be easily derived from the form of $\mathbf{X}_{(1)}$, $\mathbf{X}_{(2)}$ and $\mathbf{X}_{(3)}$.

Again, this approach allows approximate computation but comes with some difficulties, as discussed ahead.

A mosaic of uniqueness results



Inherent (or trivial) ambiguities

Typically, uniqueness of parameters of a tensor model can be only shown up to two inherent (or trivial) ambiguities:

- (i) **permutation**, since the indexing of terms in a sum is (usually) arbitrary;
- (ii) **scaling**, due to the multilinear nature of the model.

Inherent (or trivial) ambiguities

Typically, uniqueness of parameters of a tensor model can be only shown up to two inherent (or trivial) ambiguities:

- (i) **permutation**, since the indexing of terms in a sum is (usually) arbitrary;
- (ii) **scaling**, due to the multilinear nature of the model.

Example: Take $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ of size R . For any permutation $\pi \in \mathfrak{S}_R$, and any nonzero scalars α_r, β_r we have

$$\mathbf{X} = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r = \sum_{r=1}^R (\alpha_r \mathbf{a}_{\pi_r}) \otimes (\beta_r \mathbf{b}_{\pi_r}) \otimes ((\alpha_r \beta_r)^{-1} \mathbf{c}_{\pi_r}).$$

Inherent (or trivial) ambiguities

Typically, uniqueness of parameters of a tensor model can be only shown up to two **inherent (or trivial) ambiguities**:

- (i) **permutation**, since the indexing of terms in a sum is (usually) arbitrary;
- (ii) **scaling**, due to the multilinear nature of the model.

Example: Take $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ of size R . For any permutation $\pi \in \mathfrak{S}_R$, and any nonzero scalars α_r, β_r we have

$$\mathbf{X} = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r = \sum_{r=1}^R (\alpha_r \mathbf{a}_{\pi_r}) \otimes (\beta_r \mathbf{b}_{\pi_r}) \otimes ((\alpha_r \beta_r)^{-1} \mathbf{c}_{\pi_r}).$$

Partial remedy: “extract” norms λ_r and sort the terms according to them:

$$\mathbf{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r, \quad \|\mathbf{a}_r\| = \|\mathbf{b}_r\| = \|\mathbf{c}_r\| = 1, \quad \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_R \geq 0.$$

However, sign (or rotation, in \mathbb{C}) ambiguities remain.



Essential uniqueness

Def: The parameters of a tensor decomposition are essentially unique if they are unique up to trivial (permutation and scaling) ambiguities.

Essential uniqueness

Def: The parameters of a tensor decomposition are essentially unique if they are unique up to trivial (permutation and scaling) ambiguities.

Example: Essential uniqueness of a rank- R CPD $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ is equivalent to:

$$\begin{aligned} \forall (\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}) &\in \mathbb{R}^{I \times R} \times \mathbb{R}^{J \times R} \times \mathbb{R}^{K \times R}, \\ \mathbf{X} = \llbracket \tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}} \rrbracket &\Leftrightarrow \tilde{\mathbf{A}} = \mathbf{A}\mathbf{\Pi}\mathbf{\Lambda}_1, \tilde{\mathbf{B}} = \mathbf{B}\mathbf{\Pi}\mathbf{\Lambda}_2, \tilde{\mathbf{C}} = \mathbf{C}\mathbf{\Pi}\mathbf{\Lambda}_3, \end{aligned}$$

where $\mathbf{\Pi} \in \mathbb{R}^{R \times R}$ is a permutation matrix and $\mathbf{\Lambda}_i \in \mathbb{R}^{R \times R}$ are diagonal scaling matrices satisfying

$$\mathbf{\Lambda}_1 \mathbf{\Lambda}_2 \mathbf{\Lambda}_3 = \mathbf{I}.$$



Essential uniqueness

Def: The parameters of a tensor decomposition are essentially unique if they are unique up to trivial (permutation and scaling) ambiguities.

Example: Essential uniqueness of a rank- R CPD $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ is equivalent to:

$$\begin{aligned} \forall (\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}) \in \mathbb{R}^{I \times R} \times \mathbb{R}^{J \times R} \times \mathbb{R}^{K \times R}, \\ \mathbf{X} = \llbracket \tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}} \rrbracket \quad \Leftrightarrow \quad \tilde{\mathbf{A}} = \mathbf{A}\mathbf{\Pi}\mathbf{\Lambda}_1, \quad \tilde{\mathbf{B}} = \mathbf{B}\mathbf{\Pi}\mathbf{\Lambda}_2, \quad \tilde{\mathbf{C}} = \mathbf{C}\mathbf{\Pi}\mathbf{\Lambda}_3, \end{aligned}$$

where $\mathbf{\Pi} \in \mathbb{R}^{R \times R}$ is a permutation matrix and $\mathbf{\Lambda}_i \in \mathbb{R}^{R \times R}$ are diagonal scaling matrices satisfying

$$\mathbf{\Lambda}_1 \mathbf{\Lambda}_2 \mathbf{\Lambda}_3 = \mathbf{I}.$$

□

Example: Essential uniqueness of a R -block BTD

$$\mathbf{X} = \sum_{r=1}^R (\mathbf{A}_r \mathbf{B}_r)^\top \otimes \mathbf{c}_r = \sum_{r=1}^R \mathbf{H}_r^\top \otimes \mathbf{c}_r$$

means uniqueness up to permutation of blocks and to rescaling of $(\mathbf{H}_r, \mathbf{c}_r)$:

$$(\mathbf{H}_r, \mathbf{c}_r) \mapsto (\alpha \mathbf{H}_r, \alpha^{-1} \mathbf{c}_r), \quad \alpha > 0.$$

□

A first intuition

Recall: lack of uniqueness = “too much freedom”

$$\mathbf{A}\mathbf{B}^\top = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}\mathbf{B} = (\mathbf{A}\mathbf{P})(\mathbf{B}\mathbf{P}^{-1})^\top = \tilde{\mathbf{A}}\tilde{\mathbf{B}}^\top, \quad \forall \mathbf{P} \in \text{GL}_R.$$

A first intuition

Recall: lack of uniqueness = “too much freedom”

$$\mathbf{A}\mathbf{B}^\top = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}\mathbf{B} = (\mathbf{A}\mathbf{P})(\mathbf{B}\mathbf{P}^{-1})^\top = \tilde{\mathbf{A}}\tilde{\mathbf{B}}^\top, \quad \forall \mathbf{P} \in \text{GL}_R.$$

Now, take a PD of size R , say $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$, and write

$$\mathbf{X}_{(1)} = \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top.$$

If we try the same trick as above, namely

$$\mathbf{X}_{(1)} = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}(\mathbf{C} \odot \mathbf{B})^\top = \tilde{\mathbf{A}} \left((\mathbf{C} \odot \mathbf{B})\mathbf{P}^{-\top} \right)^\top,$$

then for which \mathbf{P} can we write

$$(\mathbf{C} \odot \mathbf{B})\mathbf{P}^{-\top} = \tilde{\mathbf{C}} \odot \tilde{\mathbf{B}} \quad ?$$

A first intuition

Recall: lack of uniqueness = “too much freedom”

$$\mathbf{A}\mathbf{B}^\top = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}\mathbf{B} = (\mathbf{A}\mathbf{P})(\mathbf{B}\mathbf{P}^{-1})^\top = \tilde{\mathbf{A}}\tilde{\mathbf{B}}^\top, \quad \forall \mathbf{P} \in \text{GL}_R.$$

Now, take a PD of size R , say $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$, and write

$$\mathbf{X}_{(1)} = \mathbf{A}(\mathbf{C} \odot \mathbf{B})^\top.$$

If we try the same trick as above, namely

$$\mathbf{X}_{(1)} = \mathbf{A}\mathbf{P}\mathbf{P}^{-1}(\mathbf{C} \odot \mathbf{B})^\top = \tilde{\mathbf{A}} \left((\mathbf{C} \odot \mathbf{B})\mathbf{P}^{-\top} \right)^\top,$$

then for which \mathbf{P} can we write

$$(\mathbf{C} \odot \mathbf{B})\mathbf{P}^{-\top} = \tilde{\mathbf{C}} \odot \tilde{\mathbf{B}} \quad ?$$

Rigidity: In general, $\mathbf{P}^{-\top}$ can only **permute** and **rescale** the columns of $\mathbf{C} \odot \mathbf{B}$, as these are vectorized rank-1 matrices, which do **not** form a linear space.

(Hence, \mathbf{P} only permutes and rescales the columns of \mathbf{A} as well.)

Warmup: a CPD with full rank factors

Let's prove (essential) uniqueness of the CPD when \mathbf{A} , \mathbf{B} , \mathbf{C} all have full column rank. Write

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D}, \quad \text{where} \quad \mathbf{D} = \sum_{r=1}^R \mathbf{e}_r^{\otimes 3} = \llbracket \mathbf{I}, \mathbf{I}, \mathbf{I} \rrbracket.$$

Claim: The CPD of \mathbf{X} is unique iff the CPD of \mathbf{D} is unique.

(Uniqueness of $\mathbf{D} = \llbracket \mathbf{I}, \mathbf{I}, \mathbf{I} \rrbracket$ means that every CPD of \mathbf{D} has the form $\mathbf{D} = \llbracket \mathbf{\Pi}\mathbf{\Lambda}_1, \mathbf{\Pi}\mathbf{\Lambda}_2, \mathbf{\Pi}\mathbf{\Lambda}_3 \rrbracket$.)

Warmup: a CPD with full rank factors

Let's prove (essential) uniqueness of the CPD when \mathbf{A} , \mathbf{B} , \mathbf{C} all have full column rank. Write

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D}, \quad \text{where} \quad \mathbf{D} = \sum_{r=1}^R \mathbf{e}_r^{\otimes 3} = \llbracket \mathbf{I}, \mathbf{I}, \mathbf{I} \rrbracket.$$

Claim: The CPD of \mathbf{X} is unique iff the CPD of \mathbf{D} is unique.

(Uniqueness of $\mathbf{D} = \llbracket \mathbf{I}, \mathbf{I}, \mathbf{I} \rrbracket$ means that every CPD of \mathbf{D} has the form $\mathbf{D} = \llbracket \mathbf{\Pi}\mathbf{\Lambda}_1, \mathbf{\Pi}\mathbf{\Lambda}_2, \mathbf{\Pi}\mathbf{\Lambda}_3 \rrbracket$.)

Proof: \Rightarrow) Writing a CPD $\mathbf{D} = \llbracket \mathbf{U}, \mathbf{V}, \mathbf{W} \rrbracket$ of rank R , we have

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{D} = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot [(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{D}] = \llbracket \mathbf{AU}, \mathbf{BV}, \mathbf{CW} \rrbracket.$$

By uniqueness of the CPD of \mathbf{X} , we must have

$$\mathbf{AU} = \mathbf{A}\mathbf{\Pi}\mathbf{\Lambda}_1, \quad \mathbf{BV} = \mathbf{B}\mathbf{\Pi}\mathbf{\Lambda}_2, \quad \mathbf{CW} = \mathbf{C}\mathbf{\Pi}\mathbf{\Lambda}_3$$

with $\mathbf{\Lambda}_1 \mathbf{\Lambda}_2 \mathbf{\Lambda}_3 = \mathbf{I}$. The result then follows from all factors admitting left inverses.

\Leftarrow) See exercices.



Warmup: uniqueness of $[[\mathbf{I}, \mathbf{I}, \mathbf{I}]]$

Proposition: The CPD $\mathbf{D} = [[\mathbf{I}, \mathbf{I}, \mathbf{I}]]$ is unique.

Proof sketch: Take $\mathbf{D} = [[\mathbf{U}, \mathbf{V}, \mathbf{W}]]$. From

$$\mathbf{I}(\mathbf{I} \odot \mathbf{I})^\top = \mathbf{U}(\mathbf{W} \odot \mathbf{V})^\top$$

we see that $\text{rank } \mathbf{U} = R$. By symmetry of roles, $\text{rank } \mathbf{U} = \text{rank } \mathbf{V} = \text{rank } \mathbf{W} = R$.

Warmup: uniqueness of $[[\mathbf{I}, \mathbf{I}, \mathbf{I}]]$

Proposition: The CPD $\mathbf{D} = [[\mathbf{I}, \mathbf{I}, \mathbf{I}]]$ is unique.

Proof sketch: Take $\mathbf{D} = [[\mathbf{U}, \mathbf{V}, \mathbf{W}]]$. From

$$\mathbf{I}(\mathbf{I} \odot \mathbf{I})^\top = \mathbf{U}(\mathbf{W} \odot \mathbf{V})^\top$$

we see that $\text{rank } \mathbf{U} = R$. By symmetry of roles, $\text{rank } \mathbf{U} = \text{rank } \mathbf{V} = \text{rank } \mathbf{W} = R$.

Now, write the j th slice of \mathbf{D}

$$\mathbf{e}_j \mathbf{e}_j^\top = \mathbf{U} \text{Diag}(\bar{\mathbf{v}}_j) \mathbf{W}^\top,$$

where $\bar{\mathbf{v}}_j$ is the j th row of \mathbf{V} . As \mathbf{U}, \mathbf{W} have rank R , $\bar{\mathbf{v}}_j$ can only have one nonzero entry. Suppose $\bar{\mathbf{v}}_j = c_j \mathbf{e}_{\ell_j}$, with $c_j \neq 0$. It follows that

$$\mathbf{e}_j \mathbf{e}_j^\top = c_j \mathbf{u}_{\ell_j} \mathbf{w}_{\ell_j}^\top,$$

and we are forced to take both \mathbf{u}_{ℓ_j} and \mathbf{w}_{ℓ_j} proportional to \mathbf{e}_j . Furthermore, the ℓ_j are all distinct (since $\text{rank } \mathbf{V} = R$). Hence,

$$\mathbf{U} = \mathbf{\Pi} \mathbf{\Lambda}_1, \quad \mathbf{V} = \mathbf{\Pi} \mathbf{\Lambda}_2, \quad \mathbf{W} = \mathbf{\Pi} \mathbf{\Lambda}_3,$$

with $\mathbf{\Lambda}_1 \mathbf{\Lambda}_2 \mathbf{\Lambda}_3 = \mathbf{I}$ and $\mathbf{\Pi} = (\mathbf{e}_{\ell_1} \quad \dots \quad \mathbf{e}_{\ell_R})$.



Uniqueness under relaxed assumptions

Previous result: uniqueness of CPD is already **stronger** than that of \mathbf{AB}^\top .

Uniqueness under relaxed assumptions

Previous result: uniqueness of CPD is already **stronger** than that of \mathbf{AB}^\top .

In particular, it implies:

Corollary: If $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are sampled from absolutely continuous distributions (w.r.t. Lebesgue) and $\min \{I, J, K\} \geq R$, then $[[\mathbf{A}, \mathbf{B}, \mathbf{C}]]$ is essentially unique.

Uniqueness under relaxed assumptions

Previous result: uniqueness of CPD is already **stronger** than that of \mathbf{AB}^\top .

In particular, it implies:

Corollary: If $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are sampled from absolutely continuous distributions (w.r.t. Lebesgue) and $\min \{I, J, K\} \geq R$, then $[[\mathbf{A}, \mathbf{B}, \mathbf{C}]]$ is essentially unique.

But it turns out that **uniqueness holds much more generally!**

The Kruskal rank

Def: The Kruskal rank of \mathbf{A} is the largest number $k_{\mathbf{A}}$ such that **every** subset of $k_{\mathbf{A}}$ columns of \mathbf{A} is linearly independent (= spark $\mathbf{A} - 1$ in compressive sensing).

In particular, $k_{\mathbf{A}} = 0$ iff \mathbf{A} has a null col.

Intuition: measures the “degree” of (in)dependence of the cols of $\mathbf{A} \in \mathbb{R}^{I \times R}$.

In particular,

$$k_{\mathbf{A}} \leq \text{rank } \mathbf{A}, \quad \text{with equality if } \text{rank } \mathbf{A} = R.$$

Plays a central role in Kruskal’s celebrated uniqueness result¹ & extensions.

(This is Joseph Kruskal, not to be confused with his brothers Willam Kruskal, statistician & author of the Kruskal-Wallis test, and Martin Kruskal, physicist.)

1: Kruskal, 1977

Kruskal rank: examples

In all examples below, $\{a_i\}$ forms a linearly independent family.

Example: $\mathbf{A} = (a_1 \quad a_2 \quad a_1 + a_2)$.

Kruskal rank: examples

In all examples below, $\{a_i\}$ forms a linearly independent family.

Example: $\mathbf{A} = (a_1 \quad a_2 \quad a_1 + a_2).$

$$k_{\mathbf{A}} = 2 = \text{rank } \mathbf{A} < R.$$



Example: $\mathbf{A} = (a_1 \quad a_2 \quad -2a_1).$

Kruskal rank: examples

In all examples below, $\{a_i\}$ forms a linearly independent family.

Example: $\mathbf{A} = (a_1 \quad a_2 \quad a_1 + a_2).$

$$k_{\mathbf{A}} = 2 = \text{rank } \mathbf{A} < R.$$



Example: $\mathbf{A} = (a_1 \quad a_2 \quad -2a_1).$

$k_{\mathbf{A}} = 1 < \text{rank } \mathbf{A} = 2 < R$, since \mathbf{A} holds a pair of collinear columns.



Example: $\mathbf{A} = (a_1 + a_2 \quad a_2 + a_3 \quad a_3 + a_1 \quad a_1 + 2a_2 + a_3).$

Kruskal rank: examples

In all examples below, $\{a_i\}$ forms a linearly independent family.

Example: $\mathbf{A} = (a_1 \quad a_2 \quad a_1 + a_2)$.

$$k_{\mathbf{A}} = 2 = \text{rank } \mathbf{A} < R.$$



Example: $\mathbf{A} = (a_1 \quad a_2 \quad -2a_1)$.

$$k_{\mathbf{A}} = 1 < \text{rank } \mathbf{A} = 2 < R, \text{ since } \mathbf{A} \text{ holds a pair of collinear columns.}$$



Example: $\mathbf{A} = (a_1 + a_2 \quad a_2 + a_3 \quad a_3 + a_1 \quad a_1 + 2a_2 + a_3)$.

$$k_{\mathbf{A}} = 2 < \text{rank } \mathbf{A} = 3 < R.$$



Example: $\mathbf{A} = (a_1 \quad a_2 \quad a_3 \quad 0)$.

The zero column implies $k_{\mathbf{A}} = 0$.



Some *necessary* conditions

Claim: Uniqueness of $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ cannot hold if $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}, k_{\mathbf{C}}\} < 2$.

Remark: This condition can be relaxed for higher orders.

Some *necessary* conditions

Claim: Uniqueness of $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ cannot hold if $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}, k_{\mathbf{C}}\} < 2$.

Remark: This condition can be relaxed for higher orders.

Claim: If $\text{rank } \mathbf{X} < R$, then no rank- R PD of \mathbf{X} can be unique.

Proof: Any term $\mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$ of an $(R - 1)$ -rank PD can always be split into

$$(\mathbf{a}_r - \mathbf{w}) \otimes \mathbf{b}_r \otimes \mathbf{c}_r + \mathbf{w} \otimes \mathbf{b}_r \otimes \mathbf{c}_r, \quad \forall \mathbf{w} \in \mathbb{R}^I.$$



Some *necessary* conditions

Claim: Uniqueness of $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ cannot hold if $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}, k_{\mathbf{C}}\} < 2$.

Remark: This condition can be relaxed for higher orders.

Claim: If $\text{rank } \mathbf{X} < R$, then no rank- R PD of \mathbf{X} can be unique.

Proof: Any term $\mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$ of an $(R - 1)$ -rank PD can always be split into

$$(\mathbf{a}_r - \mathbf{w}) \otimes \mathbf{b}_r \otimes \mathbf{c}_r + \mathbf{w} \otimes \mathbf{b}_r \otimes \mathbf{c}_r, \quad \forall \mathbf{w} \in \mathbb{R}^I.$$



Claim: If $\min \{\text{rank } \mathbf{B} \odot \mathbf{A}, \text{rank } \mathbf{C} \odot \mathbf{B}, \text{rank } \mathbf{C} \odot \mathbf{A}\} < R$, then $\text{rank } \mathbf{X} < R$.

Proof: If, say, $\text{rank } \mathbf{B} \odot \mathbf{A} < R$, then take any nonzero $\mathbf{v} \in \ker \mathbf{B} \odot \mathbf{A}$ and choose some \mathbf{z} such that $\mathbf{C} + \mathbf{z}\mathbf{v}^\top$ has a null column to write a rank- $(R - 1)$ PD:

$$\begin{aligned} \mathbf{X}_{(3)} &= \mathbf{C} (\mathbf{B} \odot \mathbf{A})^\top \\ &= \mathbf{C} (\mathbf{B} \odot \mathbf{A})^\top + \left[(\mathbf{B} \odot \mathbf{A}) \mathbf{v} \mathbf{z}^\top \right]^\top \\ &= (\mathbf{C} + \mathbf{z}\mathbf{v}^\top) (\mathbf{B} \odot \mathbf{A})^\top. \end{aligned}$$



Uniqueness under relaxed assumptions (cont'd)

Thm (Sidiropoulos & al., 2017): Let $\mathbf{A} \in \mathbb{R}^{I \times R}$, $\mathbf{B} \in \mathbb{R}^{J \times R}$, $\mathbf{C} \in \mathbb{R}^{R \times R}$ and $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$. If

$$\text{rank } \mathbf{C} = R \quad \text{and} \quad k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2,$$

then $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are essentially unique.

Uniqueness under relaxed assumptions (cont'd)

Thm (Sidiropoulos & al., 2017): Let $\mathbf{A} \in \mathbb{R}^{I \times R}$, $\mathbf{B} \in \mathbb{R}^{J \times R}$, $\mathbf{C} \in \mathbb{R}^{R \times R}$ and $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$. If

$$\text{rank } \mathbf{C} = R \quad \text{and} \quad k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2,$$

then $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are essentially unique.

The proof relies on a special case of [Kruskal's permutation lemma](#): the main tool used for Kruskal's 1977 proof of the most general known result.

Lemma (Sidiropoulos & al., 2017): Let $\mathbf{C}, \tilde{\mathbf{C}} \in \mathbb{R}^{R \times R}$ be two nonsingular matrices. If

$$\forall \mathbf{v} \in \mathbb{R}^K \text{ such that } \left\| \mathbf{v}^T \tilde{\mathbf{C}} \right\|_0 = 1 \quad \text{we have} \quad \left\| \mathbf{v}^T \mathbf{C} \right\|_0 = 1,$$

then $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{\Pi} \mathbf{\Lambda}_3$, where $\mathbf{\Pi}, \mathbf{\Lambda}_3$ are as before.

Uniqueness under relaxed assumptions (cont'd)

Thm (Sidiropoulos & al., 2017): Let $\mathbf{A} \in \mathbb{R}^{I \times R}$, $\mathbf{B} \in \mathbb{R}^{J \times R}$, $\mathbf{C} \in \mathbb{R}^{R \times R}$ and $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$. If

$$\text{rank } \mathbf{C} = R \quad \text{and} \quad k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2,$$

then $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are essentially unique.

The proof relies on a special case of [Kruskal's permutation lemma](#): the main tool used for Kruskal's 1977 proof of the most general known result.

Lemma (Sidiropoulos & al., 2017): Let $\mathbf{C}, \tilde{\mathbf{C}} \in \mathbb{R}^{R \times R}$ be two nonsingular matrices. If

$$\forall \mathbf{v} \in \mathbb{R}^K \text{ such that } \left\| \mathbf{v}^T \tilde{\mathbf{C}} \right\|_0 = 1 \quad \text{we have} \quad \left\| \mathbf{v}^T \mathbf{C} \right\|_0 = 1,$$

then $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{\Pi} \mathbf{\Lambda}_3$, where $\mathbf{\Pi}, \mathbf{\Lambda}_3$ are as before.

Proof: By the above condition & the fact that $\text{rank } \tilde{\mathbf{C}}^{-1} \mathbf{C} = R$,

$$\tilde{\mathbf{C}}^{-1} \tilde{\mathbf{C}} = \mathbf{I} \quad \Rightarrow \quad \tilde{\mathbf{C}}^{-1} \mathbf{C} = \mathbf{\Lambda}_3^{-1} \mathbf{\Pi}^T.$$



Uniqueness under relaxed assumptions (cont'd)

Proof (of the Thm): We will shortly see that $k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2$ implies $k_{\mathbf{B} \odot \mathbf{A}} = R = \text{rank}(\mathbf{B} \odot \mathbf{A})$. Hence, combined with

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top$$

we get $\text{rank}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = \text{rank } \tilde{\mathbf{C}} = R$.

Uniqueness under relaxed assumptions (cont'd)

Proof (of the Thm): We will shortly see that $k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2$ implies $k_{\mathbf{B} \odot \mathbf{A}} = R = \text{rank}(\mathbf{B} \odot \mathbf{A})$. Hence, combined with

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top$$

we get $\text{rank}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = \text{rank } \tilde{\mathbf{C}} = R$. Next, take a linear combination of slices

$$(\cdot, \cdot, \mathbf{v}) \cdot \mathbf{X} = \sum_{r=1}^R (\mathbf{v})_r \mathbf{X}_r = \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top = \tilde{\mathbf{A}} \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) \tilde{\mathbf{B}}^\top.$$

But, $\text{rank } \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top = \text{rank } \tilde{\mathbf{A}} \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) \tilde{\mathbf{B}}^\top \leq \text{rank } \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) = \|\mathbf{v}^\top \tilde{\mathbf{C}}\|_0$.

Suppose now $\|\mathbf{v}^\top \tilde{\mathbf{C}}\|_0 = 1$. Then $\text{rank } \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top \leq 1$ and we want to show $\|\mathbf{v}^\top \mathbf{C}\|_0 = 1$ (to use the Lemma).

Uniqueness under relaxed assumptions (cont'd)

Proof (of the Thm): We will shortly see that $k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2$ implies $k_{\mathbf{B} \odot \mathbf{A}} = R = \text{rank}(\mathbf{B} \odot \mathbf{A})$. Hence, combined with

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top$$

we get $\text{rank}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = \text{rank } \tilde{\mathbf{C}} = R$. Next, take a linear combination of slices

$$(\cdot, \cdot, \mathbf{v}) \cdot \mathbf{X} = \sum_{r=1}^R (\mathbf{v})_r \mathbf{X}_r = \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top = \tilde{\mathbf{A}} \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) \tilde{\mathbf{B}}^\top.$$

But, $\text{rank } \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top = \text{rank } \tilde{\mathbf{A}} \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) \tilde{\mathbf{B}}^\top \leq \text{rank } \text{Diag}(\mathbf{v}^\top \tilde{\mathbf{C}}) = \|\mathbf{v}^\top \tilde{\mathbf{C}}\|_0$.

Suppose now $\|\mathbf{v}^\top \tilde{\mathbf{C}}\|_0 = 1$. Then $\text{rank } \mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top \leq 1$ and we want to show $\|\mathbf{v}^\top \mathbf{C}\|_0 = 1$ (to use the Lemma).

Let $\bar{\mathbf{A}} \in \mathbb{R}^{I \times \bar{R}}$ hold the subset of $\bar{R} := \|\mathbf{v}^\top \mathbf{C}\|_0$ cols \mathbf{a}_i of \mathbf{A} such that $(\mathbf{v}^\top \mathbf{C})_i \neq 0$, similarly for $\bar{\mathbf{B}} \in \mathbb{R}^{J \times \bar{R}}$, $\bar{\mathbf{C}} \in \mathbb{R}^{R \times \bar{R}}$.

Then,

$$\mathbf{A} \text{Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top = \bar{\mathbf{A}} \text{Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top.$$

Uniqueness under relaxed assumptions (cont'd)

By the Sylvester inequality (see handout),

$$\text{rank } \bar{\mathbf{A}} \text{ Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top \geq \text{rank } \bar{\mathbf{A}} + \text{rank } \bar{\mathbf{B}} - \bar{R} \geq \min \{k_{\mathbf{A}}, \bar{R}\} + \min \{k_{\mathbf{B}}, \bar{R}\} - \bar{R}.$$

Uniqueness under relaxed assumptions (cont'd)

By the Sylvester inequality (see handout),

$$\text{rank } \bar{\mathbf{A}} \text{ Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top \geq \text{rank } \bar{\mathbf{A}} + \text{rank } \bar{\mathbf{B}} - \bar{R} \geq \min \{k_{\mathbf{A}}, \bar{R}\} + \min \{k_{\mathbf{B}}, \bar{R}\} - \bar{R}.$$

Hence, since $\text{rank } \bar{\mathbf{A}} \text{ Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top = \text{rank } \mathbf{A} \text{ Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top \leq 1$ by assumption,

$$\min \{k_{\mathbf{A}}, \bar{R}\} + \min \{k_{\mathbf{B}}, \bar{R}\} - \bar{R} \leq 1.$$

Uniqueness under relaxed assumptions (cont'd)

By the Sylvester inequality (see handout),

$$\text{rank } \bar{\mathbf{A}} \text{ Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top \geq \text{rank } \bar{\mathbf{A}} + \text{rank } \bar{\mathbf{B}} - \bar{R} \geq \min \{k_{\mathbf{A}}, \bar{R}\} + \min \{k_{\mathbf{B}}, \bar{R}\} - \bar{R}.$$

Hence, since $\text{rank } \bar{\mathbf{A}} \text{ Diag}(\mathbf{v}^\top \bar{\mathbf{C}}) \bar{\mathbf{B}}^\top = \text{rank } \mathbf{A} \text{ Diag}(\mathbf{v}^\top \mathbf{C}) \mathbf{B}^\top \leq 1$ by assumption,

$$\min \{k_{\mathbf{A}}, \bar{R}\} + \min \{k_{\mathbf{B}}, \bar{R}\} - \bar{R} \leq 1.$$

Three cases arise:

- (i) If $\bar{R} \leq \min \{k_{\mathbf{A}}, k_{\mathbf{B}}\}$, then $\bar{R} \leq 1$, hence $\bar{R} = 1$ since $\mathbf{v}^\top \mathbf{C}$ cannot vanish.
- (ii) The case $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}\} \leq \bar{R} \leq \max \{k_{\mathbf{A}}, k_{\mathbf{B}}\}$ is impossible, as it implies $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}\} \leq 1$ (violating $k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2$).
- (iii) Finally, $\max \{k_{\mathbf{A}}, k_{\mathbf{B}}\} \leq \bar{R}$ is equally impossible, since $k_{\mathbf{A}} + k_{\mathbf{B}} \geq R + 2$ by assumption, yielding

$$R + 1 \leq k_{\mathbf{A}} + k_{\mathbf{B}} - 1 \leq \bar{R} \leq R.$$

Conclusion: $\bar{R} = \|\mathbf{v}^\top \mathbf{C}\|_0 = 1$ as per the Lemma, and therefore $\tilde{\mathbf{C}} = \mathbf{C} \Pi \Lambda_3$.

Uniqueness under relaxed assumptions (cont'd)

Now, from

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top = \mathbf{C}\Pi\Lambda_3(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top,$$

we pre-multiply by \mathbf{C}^{-1} to get

$$(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = (\mathbf{B} \odot \mathbf{A})\Pi\Lambda_3^{-1} = (\mathbf{B}\Pi\Lambda_2) \odot (\mathbf{A}\Pi\Lambda_1),$$

where $\Lambda_1\Lambda_2 = \Lambda_3^{-1}$, as claimed. ■

Uniqueness under relaxed assumptions (cont'd)

Now, from

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top = \mathbf{C}\Pi\Lambda_3(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top,$$

we pre-multiply by \mathbf{C}^{-1} to get

$$(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = (\mathbf{B} \odot \mathbf{A})\Pi\Lambda_3^{-1} = (\mathbf{B}\Pi\Lambda_2) \odot (\mathbf{A}\Pi\Lambda_1),$$

where $\Lambda_1\Lambda_2 = \Lambda_3^{-1}$, as claimed. ■

Remark: The extension to tall $\mathbf{C} \in \mathbb{R}^{K \times R}$ is easy: if $\text{rank } \mathbf{C} = R$, then it contains a nonsingular $R \times R$ submatrix. Uniqueness and full rank of $\mathbf{A} \odot \mathbf{B}$ then yields uniqueness of the whole \mathbf{C} :

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top = \tilde{\mathbf{C}} [(\mathbf{B} \odot \mathbf{A})\Pi\Lambda_3^{-1}]^\top$$

Uniqueness under relaxed assumptions (cont'd)

Now, from

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top = \mathbf{C}\Pi\Lambda_3(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top,$$

we pre-multiply by \mathbf{C}^{-1} to get

$$(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}}) = (\mathbf{B} \odot \mathbf{A})\Pi\Lambda_3^{-1} = (\mathbf{B}\Pi\Lambda_2) \odot (\mathbf{A}\Pi\Lambda_1),$$

where $\Lambda_1\Lambda_2 = \Lambda_3^{-1}$, as claimed. ■

Remark: The extension to tall $\mathbf{C} \in \mathbb{R}^{K \times R}$ is easy: if $\text{rank } \mathbf{C} = R$, then it contains a nonsingular $R \times R$ submatrix. Uniqueness and full rank of $\mathbf{A} \odot \mathbf{B}$ then yields uniqueness of the whole \mathbf{C} :

$$\mathbf{C}(\mathbf{B} \odot \mathbf{A})^\top = \tilde{\mathbf{C}}(\tilde{\mathbf{B}} \odot \tilde{\mathbf{A}})^\top = \tilde{\mathbf{C}} [(\mathbf{B} \odot \mathbf{A})\Pi\Lambda_3^{-1}]^\top$$

Example: If $R = 3$ and

$$\mathbf{A} = (a_1 \quad a_2 \quad a_3), \quad \mathbf{B} = (b_1 \quad b_2 \quad b_1 + b_2), \quad \mathbf{C} = (c_1 \quad c_2 \quad c_3),$$

where distinct vectors are independent, then $\text{rank } \mathbf{C} = 3 = R$ and $k_{\mathbf{A}} + k_{\mathbf{B}} = 5 = R + 2$, hence $[\mathbf{A}, \mathbf{B}, \mathbf{C}]$ is unique. □

Kruskal's uniqueness theorem

Jiang & Sidiropoulos (2004) gave another, more accessible proof of Kruskal's permutation lemma in 2004. A simplified statement is:

Lemma (Stegeman & Sidiropoulos, 2007): Let $\mathbf{C}, \tilde{\mathbf{C}} \in \mathbb{R}^{K \times R}$ and $k_{\mathbf{C}} \geq 2$. If

$$\forall \mathbf{v} \in \mathbb{R}^K, \quad \left\| \mathbf{v}^T \mathbf{C} \right\|_0 \leq \left\| \mathbf{v}^T \tilde{\mathbf{C}} \right\|_0,$$

then $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{\Pi} \mathbf{\Lambda}_3$, where $\mathbf{\Pi}, \mathbf{\Lambda}_3$ are as before.

Kruskal's uniqueness theorem

Jiang & Sidiropoulos (2004) gave another, more accessible proof of Kruskal's permutation lemma in 2004. A simplified statement is:

Lemma (Stegeman & Sidiropoulos, 2007): Let $\mathbf{C}, \tilde{\mathbf{C}} \in \mathbb{R}^{K \times R}$ and $k_{\mathbf{C}} \geq 2$. If

$$\forall \mathbf{v} \in \mathbb{R}^K, \quad \left\| \mathbf{v}^T \mathbf{C} \right\|_0 \leq \left\| \mathbf{v}^T \tilde{\mathbf{C}} \right\|_0,$$

then $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{\Pi} \mathbf{\Lambda}_3$, where $\mathbf{\Pi}, \mathbf{\Lambda}_3$ are as before.

This is the main tool used in the proof of Kruskal's famous uniqueness result:

Thm (Kruskal, 1977): Let $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$ of size $R > 1$, such that

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} \geq 2R + 2.$$

Then, $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are essentially unique (and $\text{rank } \mathbf{X} = R$).

Kruskal's uniqueness theorem

Jiang & Sidiropoulos (2004) gave another, more accessible proof of Kruskal's permutation lemma in 2004. A simplified statement is:

Lemma (Stegeman & Sidiropoulos, 2007): Let $\mathbf{C}, \tilde{\mathbf{C}} \in \mathbb{R}^{K \times R}$ and $k_{\mathbf{C}} \geq 2$. If

$$\forall \mathbf{v} \in \mathbb{R}^K, \quad \left\| \mathbf{v}^\top \mathbf{C} \right\|_0 \leq \left\| \mathbf{v}^\top \tilde{\mathbf{C}} \right\|_0,$$

then $\tilde{\mathbf{C}} = \mathbf{C} \mathbf{\Pi} \mathbf{\Lambda}_3$, where $\mathbf{\Pi}, \mathbf{\Lambda}_3$ are as before.

This is the main tool used in the proof of Kruskal's famous uniqueness result:

Thm (Kruskal, 1977): Let $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$ of size $R > 1$, such that

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} \geq 2R + 2.$$

Then, $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are essentially unique (and $\text{rank } \mathbf{X} = R$).

Generically, $k_{\mathbf{A}} = \min \{I, R\}$, $k_{\mathbf{B}} = \min \{J, R\}$ and $k_{\mathbf{C}} = \min \{K, R\}$, so that uniqueness holds whenever $R \leq \frac{1}{2}(\min \{I, R\} + \min \{J, R\} + \min \{K, R\} - 2)$.

Is Kruskal's condition also necessary?

For $R = 2$, the condition $k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} \geq 2R + 2 = 6$ is also **necessary**, since otherwise $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}, k_{\mathbf{C}}\} < 2$.

Is Kruskal's condition also necessary?

For $R = 2$, the condition $k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} \geq 2R + 2 = 6$ is also **necessary**, since otherwise $\min \{k_{\mathbf{A}}, k_{\mathbf{B}}, k_{\mathbf{C}}\} < 2$.

ten Berge and Sidiropoulos (2002) examined necessity for higher R :

- Still necessary for $R = 3$ (proof by enumeration of cases).
- Necessary also for $R = 4$ **if ranks = k-ranks**. Otherwise, no (see exercises).
- Not necessary in general for $R > 5$. Finding necessary and sufficient conditions is complicated:

“It also has been shown that, in cases of small k-rank, the particular pattern of zeros, after pretransformation to have identity submatrices in \mathbf{A} , \mathbf{B} , \mathbf{C} , may have a decisive impact on uniqueness. This implies that attempts to derive necessary and sufficient conditions for uniqueness are doomed unless they take that very pattern into account.”

Extension to higher orders

Kruskal's theorem was extended by Sidiropoulos & Bro (2000) by reduction to the third-order case ($d = 3$) via the following additivity lemma:

Lemma (Sidiropoulos & Bro, 2000): Let $\mathbf{A} \in \mathbb{R}^{I \times R}$, $\mathbf{B} \in \mathbb{R}^{J \times R}$. If $k_{\mathbf{A}} k_{\mathbf{B}} \neq 0$, then

$$R \geq k_{\mathbf{A} \odot \mathbf{B}} \geq \min \{R, k_{\mathbf{A}} + k_{\mathbf{B}} - 1\}.$$

(Otherwise, $k_{\mathbf{A} \odot \mathbf{B}} = 0$.)

Extension to higher orders

Kruskal's theorem was extended by Sidiropoulos & Bro (2000) by reduction to the third-order case ($d = 3$) via the following additivity lemma:

Lemma (Sidiropoulos & Bro, 2000): Let $\mathbf{A} \in \mathbb{R}^{I \times R}$, $\mathbf{B} \in \mathbb{R}^{J \times R}$. If $k_{\mathbf{A}} k_{\mathbf{B}} \neq 0$, then

$$R \geq k_{\mathbf{A} \odot \mathbf{B}} \geq \min \{R, k_{\mathbf{A}} + k_{\mathbf{B}} - 1\}.$$

(Otherwise, $k_{\mathbf{A} \odot \mathbf{B}} = 0$.)

For instance, for $d = 4$ and $\mathbf{X} = [\mathbf{A}, \mathbf{B}, \mathbf{G}, \mathbf{H}]_R$:

$$x_{ijkl} = \sum_{r=1}^R (\mathbf{a}_r)_i (\mathbf{b}_r)_j (\mathbf{g}_r)_k (\mathbf{h}_r)_\ell, \quad \text{with } k_{\mathbf{A}} \geq k_{\mathbf{B}} \geq k_{\mathbf{G}} \geq k_{\mathbf{H}},$$

define $\mathbf{Y} \in \mathbb{R}^{I \times J \times M}$ with $M = KL$ via the bijection $(k, \ell) \leftrightarrow m$:

$$y_{ijm(k,\ell)} = x_{ijkl}, \quad m(k, \ell) = (k - 1)L + \ell$$

Extension to higher orders (cont'd)

This implies

$$y_{ijm} = \sum_{r=1}^R (\mathbf{a}_r)_i (\mathbf{b}_r)_j (\mathbf{c}_r)_m, \quad \text{with} \quad \mathbf{C} = \mathbf{G} \odot \mathbf{H}.$$

Extension to higher orders (cont'd)

This implies

$$y_{ijm} = \sum_{r=1}^R (\mathbf{a}_r)_i (\mathbf{b}_r)_j (\mathbf{c}_r)_m, \quad \text{with} \quad \mathbf{C} = \mathbf{G} \odot \mathbf{H}.$$

If this decomposition is unique, so is that of \mathbf{X} . By Kruskal's thm, this holds when

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} = k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{G} \odot \mathbf{H}} \geq 2R + 2.$$

We have two cases:

(i) If $k_{\mathbf{G}} + k_{\mathbf{H}} - 1 \leq R$, then by the lemma $k_{\mathbf{C}} \geq k_{\mathbf{G}} + k_{\mathbf{H}} - 1$. In that case, if

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{G}} + k_{\mathbf{H}} \geq 2R + 3,$$

then Kruskal's bound is met.

Extension to higher orders (cont'd)

This implies

$$y_{ijm} = \sum_{r=1}^R (\mathbf{a}_r)_i (\mathbf{b}_r)_j (\mathbf{c}_r)_m, \quad \text{with} \quad \mathbf{C} = \mathbf{G} \odot \mathbf{H}.$$

If this decomposition is unique, so is that of \mathbf{X} . By Kruskal's thm, this holds when

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} = k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{G} \odot \mathbf{H}} \geq 2R + 2.$$

We have two cases:

(i) If $k_{\mathbf{G}} + k_{\mathbf{H}} - 1 \leq R$, then by the lemma $k_{\mathbf{C}} \geq k_{\mathbf{G}} + k_{\mathbf{H}} - 1$. In that case, if

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{G}} + k_{\mathbf{H}} \geq 2R + 3,$$

then Kruskal's bound is met.

(ii) If $k_{\mathbf{G}} + k_{\mathbf{H}} - 1 > R$, then $k_{\mathbf{C}} = R$. Combined with $k_{\mathbf{A}} + k_{\mathbf{B}} \geq k_{\mathbf{G}} + k_{\mathbf{H}} > R + 1$, we get

$$k_{\mathbf{A}} + k_{\mathbf{B}} + k_{\mathbf{C}} \geq 2R + 2.$$

Extension to higher orders (cont'd)

This argument can be generalized, yielding:

Thm (Sidiropoulos & Bro, 2000): If

$$\sum_{i=1}^d k_{\mathbf{A}^{(i)}} \geq 2R + d - 1,$$

then the decomposition $\mathbf{X} = [\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(d)}]_R$ is essentially unique.

Extension to higher orders (cont'd)

This argument can be generalized, yielding:

Thm (Sidiropoulos & Bro, 2000): If

$$\sum_{i=1}^d k_{\mathbf{A}^{(i)}} \geq 2R + d - 1,$$

then the decomposition $\mathbf{X} = [\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(d)}]_R$ is essentially unique.

Remarks:

- Uniqueness now does not need $\min_i k_{\mathbf{A}^{(i)}} > 1$. Example: $R = 2, d = 4$, $k_{\mathbf{A}^{(1)}} = k_{\mathbf{A}^{(2)}} = k_{\mathbf{A}^{(3)}} = R = 2$ and $k_{\mathbf{A}^{(4)}} = 1$.
- The condition becomes less restrictive as d grows (since it “distributes” the term $2R - 1$ among d factors).

Generic uniqueness

Generic factors: If the columns of each $\mathbf{A}^{(i)} \in \mathbb{R}^{N_i \times R}$ are **independently** drawn from an abs. continuous distribution, then $\mathbb{P}(k_{\mathbf{A}^{(i)}} = \min \{N_i, R\}) = 1$.

In this case, the sufficient condition reduces to

$$\sum_{i=1}^d \min \{N_i, R\} \geq 2R + d - 1.$$

In particular, if d grows by one, the LHS grows by at least two (since $N_i \geq 2$). Hence, the bound is eventually satisfied as d grows.

Generic uniqueness

Generic factors: If the columns of each $\mathbf{A}^{(i)} \in \mathbb{R}^{N_i \times R}$ are **independently** drawn from an abs. continuous distribution, then $\mathbb{P}(k_{\mathbf{A}^{(i)}} = \min \{N_i, R\}) = 1$.

In this case, the sufficient condition reduces to

$$\sum_{i=1}^d \min \{N_i, R\} \geq 2R + d - 1.$$

In particular, if d grows by one, the LHS grows by at least two (since $N_i \geq 2$). Hence, the bound is eventually satisfied as d grows.

Generic tensor of subgeneric rank: A generic tensor with rank bounded as

$$R \leq \text{grank}(N_1, \dots, N_d) - 1 = \left\lceil \frac{\prod_{i=1}^d N_i}{1 + \sum_{i=1}^d (N_i - 1)} \right\rceil - 1$$

admits a unique CPD almost surely if $\prod_{i=1}^d N_i \leq 15000$ (w/ some exceptions).¹

(Compare the above bounds on R .)

1: Chiantini & al., 2014

Uniqueness of the block-term decomposition

Model:

$$X = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{c}_r = \sum_{r=1}^R (\overset{N_1 \times L_r}{\mathbf{A}_r} \overset{N_2 \times L_r}{\mathbf{B}_r^\top}) \otimes \mathbf{c}_r, \quad \text{rank } \mathbf{H}_r \leq L_r$$

Factors: $\mathbf{A} = (\mathbf{A}_1 \ \dots \ \mathbf{A}_R)$, $\mathbf{B} = (\mathbf{B}_1 \ \dots \ \mathbf{B}_R)$ and $\mathbf{C} = (\mathbf{c}_1 \ \dots \ \mathbf{c}_R)$.

Uniqueness of the block-term decomposition

Model:

$$X = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{c}_r = \sum_{r=1}^R (\overset{N_1 \times L_r}{\mathbf{A}_r} \overset{N_2 \times L_r}{\mathbf{B}_r^T}) \otimes \mathbf{c}_r, \quad \text{rank } \mathbf{H}_r \leq L_r$$

Factors: $\mathbf{A} = (\mathbf{A}_1 \ \dots \ \mathbf{A}_R)$, $\mathbf{B} = (\mathbf{B}_1 \ \dots \ \mathbf{B}_R)$ and $\mathbf{C} = (\mathbf{c}_1 \ \dots \ \mathbf{c}_R)$.

Thm (De Lathauwer, 2008): If $k_{\mathbf{C}} > 1$ and \mathbf{A}, \mathbf{B} are full column rank, then the above BTM essentially unique.

Remark: This requires (and holds generically when) $\sum_r L_r \leq \min\{N_1, N_2\}$.

Uniqueness of the block-term decomposition

Model:

$$\mathbf{X} = \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{c}_r = \sum_{r=1}^R (\overset{N_1 \times L_r}{\mathbf{A}_r} \overset{N_2 \times L_r}{\mathbf{B}_r^T}) \otimes \mathbf{c}_r, \quad \text{rank } \mathbf{H}_r \leq L_r$$

Factors: $\mathbf{A} = (\mathbf{A}_1 \ \dots \ \mathbf{A}_R)$, $\mathbf{B} = (\mathbf{B}_1 \ \dots \ \mathbf{B}_R)$ and $\mathbf{C} = (\mathbf{c}_1 \ \dots \ \mathbf{c}_R)$.

Thm (De Lathauwer, 2008): If $k_{\mathbf{C}} > 1$ and \mathbf{A}, \mathbf{B} are full column rank, then the above BTM essentially unique.

Remark: This requires (and holds generically when) $\sum_r L_r \leq \min\{N_1, N_2\}$.

Thm (De Lathauwer, 2011): If \mathbf{C} has full column rank and every nontrivial combination of \mathbf{H}_r yields a matrix of higher rank:

$$\text{rank} \sum_{\ell=1}^p w_{\ell} \mathbf{H}_{r_{\ell}} > \max_{\ell} \text{rank } \mathbf{H}_{r_{\ell}}, \quad w_{\ell} \neq 0,$$

then the BTM is unique. The latter condition is necessary.

Example: $\mathbf{H}_1 \otimes \mathbf{c}_1 + \mathbf{H}_2 \otimes \mathbf{c}_2 = (\mathbf{H}_1 + w\mathbf{H}_2) \otimes \mathbf{c}_1 + \mathbf{H}_2 \otimes (\mathbf{c}_2 - w\mathbf{c}_1).$

□

Rank-1 approx., tensor spectrum & power iteration



Low-rank decomposition

Up to now: focus on the properties of various (exact) tensor decompositions.

Hereafter: how to compute such decompositions?

Low-rank decomposition

Up to now: focus on the properties of various (exact) tensor decompositions.

Hereafter: how to compute such decompositions?

Warmup: Recovering the vectors that make up a rank-1 tensor

$$\mathbf{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

Low-rank decomposition

Up to now: focus on the properties of various (exact) tensor decompositions.

Hereafter: how to compute such decompositions?

Warmup: Recovering the vectors that make up a rank-1 tensor

$$\mathbf{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

That's easy: compute the dominant left singular vectors of

$$\mathbf{X}_{(1)} = \mathbf{a}(\mathbf{c} \boxtimes \mathbf{b})^\top, \quad \mathbf{X}_{(2)} = \mathbf{b}(\mathbf{c} \boxtimes \mathbf{a})^\top, \quad \mathbf{X}_{(3)} = \mathbf{c}(\mathbf{b} \boxtimes \mathbf{a})^\top$$

Low-rank decomposition

Up to now: focus on the properties of various (exact) tensor decompositions.

Hereafter: how to compute such decompositions?

Warmup: Recovering the vectors that make up a rank-1 tensor

$$\mathbf{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$$

That's easy: compute the dominant left singular vectors of

$$\mathbf{X}_{(1)} = \mathbf{a}(\mathbf{c} \boxtimes \mathbf{b})^\top, \quad \mathbf{X}_{(2)} = \mathbf{b}(\mathbf{c} \boxtimes \mathbf{a})^\top, \quad \mathbf{X}_{(3)} = \mathbf{c}(\mathbf{b} \boxtimes \mathbf{a})^\top$$

Now let's make things more interesting: recover $\mathbf{a}, \mathbf{b}, \mathbf{c}$ from

$$\mathbf{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} + \mathbf{W},$$

where \mathbf{W} is a noise term ($\mathbb{R}^{N_1 \times N_2 \times N_3}$ -valued realization of a random tensor).

What about using the same solution?

Low-rank approximation (LRA)

Computing the dominant singular vectors of X now amounts to looking for a **rank-one approximation** of it.

It is a natural way of estimating the low-rank signal $a \otimes b \otimes c$ “planted” in X (though not optimal, as we’ll see).

More generally, low-rank models are virtually always computed by means of **low-rank approximation (LRA)** algorithms.

Low-rank approximation (LRA)

Computing the dominant singular vectors of X now amounts to looking for a **rank-one approximation** of it.

It is a natural way of estimating the low-rank signal $a \otimes b \otimes c$ “planted” in X (though not optimal, as we’ll see).

More generally, low-rank models are virtually always computed by means of **low-rank approximation (LRA)** algorithms.

Several natural questions arise:

- In which sense should one approximate X ?
- How should one formulate the approximation problem?
- Which algorithms can be deployed, and how do they perform?
- What is the best possible performance among all possible algorithms?
- Under which conditions can LRA recover the sought signal/information?

A rank-1 model with Gaussian noise

Take the rank-1-Gaussian model (symmetric, for simplicity) of size $N \times \cdots \times N$:

$$\mathbf{X} = \lambda \mathbf{a}^{\otimes 3} + \mathbf{W}, \quad \mathbf{a}^{\otimes 3} := \mathbf{a} \otimes \mathbf{a} \otimes \mathbf{a},$$

where $\|\mathbf{a}\| = 1$, $\lambda \in \mathbb{R}$ is an SNR parameter and \mathbf{W} is Gaussian & symmetric:

$$p(\mathbf{W}) = \frac{1}{Z_3(N)} \exp \left(-\frac{N}{2} \|\mathbf{W}\|_{\text{F}}^2 \right)$$

A rank-1 model with Gaussian noise

Take the rank-1-Gaussian model (symmetric, for simplicity) of size $N \times \cdots \times N$:

$$\mathbf{X} = \lambda \mathbf{a}^{\otimes 3} + \mathbf{W}, \quad \mathbf{a}^{\otimes 3} := \mathbf{a} \otimes \mathbf{a} \otimes \mathbf{a},$$

where $\|\mathbf{a}\| = 1$, $\lambda \in \mathbb{R}$ is an SNR parameter and \mathbf{W} is Gaussian & symmetric:

$$p(\mathbf{W}) = \frac{1}{Z_3(N)} \exp \left(-\frac{N}{2} \|\mathbf{W}\|_{\text{F}}^2 \right)$$

Consequently, $p(\mathbf{X} | \lambda, \mathbf{a}) \sim \exp \left(-\frac{N}{2} \|\mathbf{X} - \lambda \mathbf{a}^{\otimes 3}\|_{\text{F}}^2 \right)$ and:

Maximum likelihood estimator (MLE):

$$(\hat{\lambda}, \hat{\mathbf{a}}) := \arg \min_{\mu, \|\mathbf{u}\|=1} \|\mathbf{X} - \mu \mathbf{u}^{\otimes 3}\|_{\text{F}}^2$$

Best rank-1 approx. & spectral tensor norm

Well-known equivalence:

$$\|X\| := \max_{\|u\|=1} \left| \sum_{ijk} x_{ijk} u_i u_j u_k \right| \Leftrightarrow \min_{\mu, \|u\|=1} \overbrace{\sum_{ijk} (x_{ijk} - \mu u_i u_j u_k)^2}^{\|X - \mu u^{\otimes 3}\|_F^2}$$

Best rank-1 approx. & spectral tensor norm

Well-known equivalence:

$$\|X\| := \max_{\|u\|=1} \left| \sum_{ijk} x_{ijk} u_i u_j u_k \right| \Leftrightarrow \min_{\mu, \|u\|=1} \overbrace{\sum_{ijk} (x_{ijk} - \mu u_i u_j u_k)^2}^{\|X - \mu u^{\otimes 3}\|_F^2}$$

Proof:

$$\begin{aligned} F(\mu, u) &= \|X - \mu u^{\otimes 3}\|_F^2 = \|X\|_F^2 - 2\mu \underbrace{\langle X, u^{\otimes 3} \rangle}_{=(u, u, u) \cdot X} + \mu^2 \underbrace{\|u^{\otimes 3}\|_F^2}_{=1} \\ &= \|X\|_F^2 + \mu (\mu - 2 (u, u, u) \cdot X) \end{aligned}$$

Setting $\frac{\partial F}{\partial \mu}(\mu, u) = 0$ gives $\mu = (u, u, u) \cdot X$, leading to:

$$\min_{\|u\|=1} -((u, u, u) \cdot X)^2 \Leftrightarrow \max_{\|u\|=1} \left| (u, u, u) \cdot X \right|$$



Best rank-1 approx. & spectral tensor norm

Well-known equivalence:

$$\|X\| := \max_{\|u\|=1} \left| \sum_{ijk} x_{ijk} u_i u_j u_k \right| \Leftrightarrow \min_{\mu, \|u\|=1} \overbrace{\sum_{ijk} (x_{ijk} - \mu u_i u_j u_k)^2}^{\|X - \mu u^{\otimes 3}\|_F^2}$$

Proof:

$$\begin{aligned} F(\mu, u) &= \|X - \mu u^{\otimes 3}\|_F^2 = \|X\|_F^2 - 2\mu \underbrace{\langle X, u^{\otimes 3} \rangle}_{=(u, u, u) \cdot X} + \mu^2 \underbrace{\|u^{\otimes 3}\|_F^2}_{=1} \\ &= \|X\|_F^2 + \mu (\mu - 2(u, u, u) \cdot X) \end{aligned}$$

Setting $\frac{\partial F}{\partial \mu}(\mu, u) = 0$ gives $\mu = (u, u, u) \cdot X$, leading to:

$$\min_{\|u\|=1} -((u, u, u) \cdot X)^2 \Leftrightarrow \max_{\|u\|=1} |(u, u, u) \cdot X|$$

$$\begin{aligned} (-u, -u, -u) \cdot X &= -(u, u, u) \cdot X \\ &\vdots \end{aligned}$$

Moreover: $\max_{\|u\|=1} |(u, u, u) \cdot X| \Leftrightarrow \max_{\|u\|=1} (u, u, u) \cdot X$



Spectral norm & tensor eigenpairs

For convenience, denote $X \cdot a^3 := (a, a, a) \cdot X$ and $X \cdot a^2 := (a, a, \cdot) \cdot X$.

MLE problem

$$\max_{\|u\|=1} X \cdot u^3$$

Lagrangian

$$L(\mu, u) = \frac{1}{3} X \cdot u^3 - \frac{\mu}{2} (\|u\|^2 - 1)$$

Critical points satisfy: $\frac{\partial}{\partial u} L(\mu, u) = X \cdot u^2 - \mu u = 0, \quad \|u\| = 1$

Spectral norm & tensor eigenpairs

For convenience, denote $\mathbf{X} \cdot \mathbf{a}^3 := (\mathbf{a}, \mathbf{a}, \mathbf{a}) \cdot \mathbf{X}$ and $\mathbf{X} \cdot \mathbf{a}^2 := (\mathbf{a}, \mathbf{a}, \cdot) \cdot \mathbf{X}$.

MLE problem

$$\max_{\|\mathbf{u}\|=1} \mathbf{X} \cdot \mathbf{u}^3$$

Lagrangian

$$L(\mu, \mathbf{u}) = \frac{1}{3} \mathbf{X} \cdot \mathbf{u}^3 - \frac{\mu}{2} (\|\mathbf{u}\|^2 - 1)$$

Critical points satisfy: $\frac{\partial}{\partial \mathbf{u}} L(\mu, \mathbf{u}) = \mathbf{X} \cdot \mathbf{u}^2 - \mu \mathbf{u} = 0, \quad \|\mathbf{u}\| = 1$

Def: Tensor ℓ_2 -eigenvalue equations¹: $\mathbf{X} \cdot \mathbf{u}^2 = \mu \mathbf{u}, \quad \|\mathbf{u}\| = 1$

In particular, the MLE $\hat{\mathbf{a}}$ (any global maximizer) verifies $\mathbf{X} \cdot \hat{\mathbf{a}}^2 = \mu_{\max} \hat{\mathbf{a}}$.

But how many critical points should one expect to exist?

1: Lim, 2005

Spectral norm & tensor eigenpairs

For convenience, denote $X \cdot a^3 := (a, a, a) \cdot X$ and $X \cdot a^2 := (a, a, \cdot) \cdot X$.

MLE problem

$$\max_{\|u\|=1} X \cdot u^3$$

Lagrangian

$$L(\mu, u) = \frac{1}{3} X \cdot u^3 - \frac{\mu}{2} (\|u\|^2 - 1)$$

Critical points satisfy: $\frac{\partial}{\partial u} L(\mu, u) = X \cdot u^2 - \mu u = 0, \quad \|u\| = 1$

Def: Tensor ℓ_2 -eigenvalue equations¹: $X \cdot u^2 = \mu u, \quad \|u\| = 1$

In particular, the MLE \hat{a} (any global maximizer) verifies $X \cdot \hat{a}^2 = \mu_{\max} \hat{a}$.

But how many critical points should one expect to exist?

Thm (Cartwright & Sturmfels, 2013): A symmetric tensor of dims N has up to $((d-1)^N - 1)/(d-2)$ distinct eigenvalues (the bound is generically attained).

1: Lim, 2005

Algorithm: tensor power iteration

The tensor eigenvalue eqns naturally suggest a fixed point iteration:

$$\hat{\mathbf{u}} \leftarrow (\mathbf{X} \cdot \hat{\mathbf{u}}^{d-1}) / \hat{\mu}, \quad \text{with} \quad \hat{\mu} = \left\| \mathbf{X} \cdot \hat{\mathbf{u}}^{d-1} \right\|$$

This is known as tensor power iteration.^{1,2}

1: De Lathauwer & al., 2000b, 2: Kofidis & Regaila, 2002, 3: Robeva, 2016,
4: Kolda, 2001

Algorithm: tensor power iteration

The tensor eigenvalue eqns naturally suggest a fixed point iteration:

$$\hat{\mathbf{u}} \leftarrow (\mathbf{X} \cdot \hat{\mathbf{u}}^{d-1}) / \hat{\mu}, \quad \text{with} \quad \hat{\mu} = \left\| \mathbf{X} \cdot \hat{\mathbf{u}}^{d-1} \right\|$$

This is known as tensor power iteration.^{1,2}

It succeeds (with proper initialization) in finding all rank-1 terms of an orthogonal decomposition (**odeco**^{3,4}) $\mathbf{X} = \sum_{r=1}^R \lambda_r \mathbf{v}_r^{\otimes d}$, where $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}$:

Thm (Anandkumar & al., 2014): Let $\mathbf{X} = \sum_{r=1}^R \lambda_r \mathbf{v}_r^{\otimes d}$ be an odeco of \mathbf{X} . Then, with random initialization, power iteration converges to one \mathbf{v}_r a.s.

Furthermore, $\{\mathbf{v}_r\}$ are the only **robust eigenvectors** of \mathbf{X} : those having a neighborhood where power iteration converges to them.

But what about the non-orthogonal case?

1: De Lathauwer & al., 2000b, 2: Kofidis & Regaila, 2002, 3: Robeva, 2016,
4: Kolda, 2001

Shifted tensor power iteration

Most symmetric tensors are **not** odeco (generically, they aren't).

Still, power iteration is useful for rank-1 approximation (ex: tensor PCA).

However, the basic iteration (previous slide) is not generally convergent.¹

1: Kolda & Mayo, 2011

Shifted tensor power iteration

Most symmetric tensors are **not** odeco (generically, they aren't).

Still, power iteration is useful for rank-1 approximation (ex: tensor PCA).

However, the basic iteration (previous slide) is not generally convergent.¹

Kolda & Mayo (2011) have thus proposed a shifted power iteration:

$$\hat{\mathbf{u}} \leftarrow (\mathbf{X} \cdot \hat{\mathbf{u}}^{d-1} + \alpha \hat{\mathbf{u}}) / \hat{\mu}, \quad \text{with} \quad \hat{\mu} = \left\| \mathbf{X} \cdot \hat{\mathbf{u}}^{d-1} + \alpha \hat{\mathbf{u}} \right\|$$

They established the following sufficient condition for convergence:

$$|\alpha| > (d-1) \max_{\|\mathbf{v}\|=1} \left\| \mathbf{X} \cdot \mathbf{v}^{d-2} \right\|.$$

This allows computing all **stable** eigenpairs (μ, \mathbf{u}) : those at which $\nabla_{\mathbf{u}\mathbf{u}}^2 L(\mu, \mathbf{u})$ is definite when projected onto \mathbf{u}^\perp . (By contrast, only the dominant eigenpair can be computed in the matrix case.)

¹: Kolda & Mayo, 2011

Low-rank approximation, in several flavors



Low-rank tensor approximation

We turn now to **low-rank approximation** of a tensor.

In the matrix case, a central result is:

Thm (Eckart–Young, 1936): The solution to

$$\min_{\text{rank } \hat{\mathbf{X}} \leq R} \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F$$

corresponds to truncating the SVD of \mathbf{X} at rank R .

Low-rank tensor approximation

We turn now to **low-rank approximation** of a tensor.

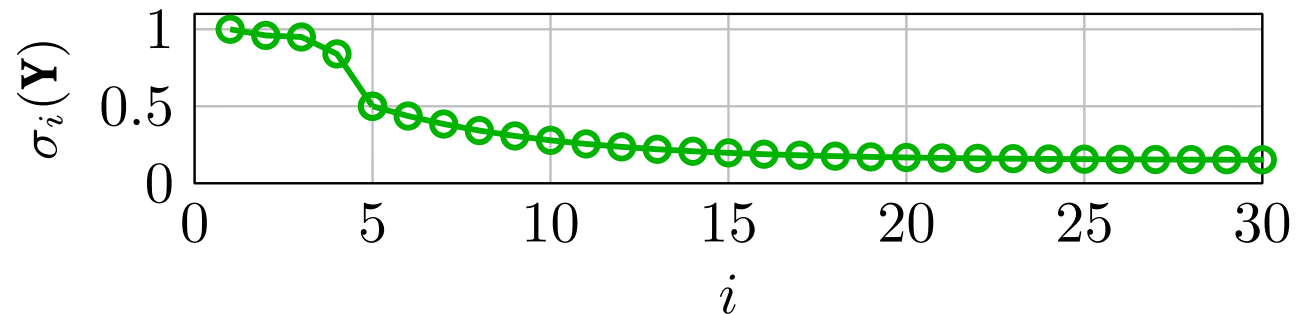
In the matrix case, a central result is:

Thm (Eckart–Young, 1936): The solution to

$$\min_{\text{rank } \hat{\mathbf{X}} \leq R} \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F$$

corresponds to truncating the SVD of \mathbf{X} at rank R .

$$\mathbf{Y} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$



Low-rank tensor approximation

We turn now to **low-rank approximation** of a tensor.

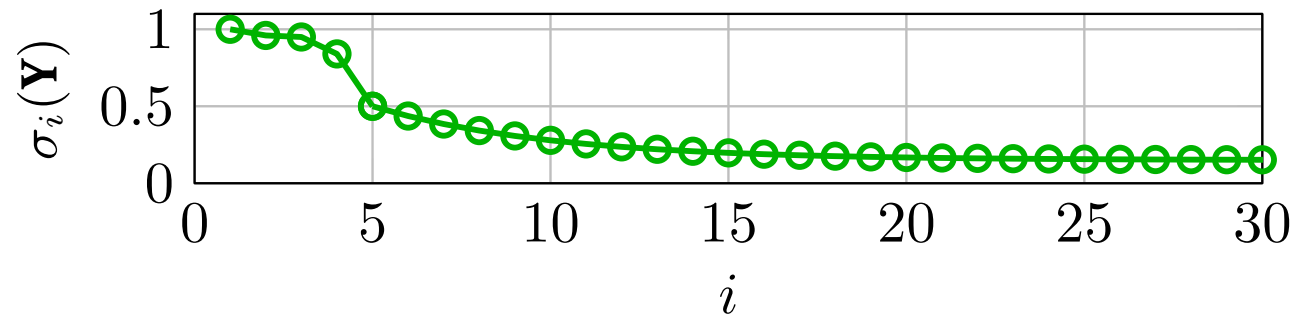
In the matrix case, a central result is:

Thm (Eckart–Young, 1936): The solution to

$$\min_{\text{rank } \hat{\mathbf{X}} \leq R} \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F$$

corresponds to truncating the SVD of \mathbf{X} at rank R .

$$\mathbf{Y} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$



No such a general result exists for tensors, for the usual notions of rank.

Moreover, for rank > 1 , all the discussion on results & algorithms is conditioned upon the definition of rank which is relevant to a certain end.

Low-mrank approximation

Goal: Approximate $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$ by $\hat{\mathbf{X}}$ such that $\text{mrank } \hat{\mathbf{X}} \leq (R_1, R_2, R_3)$

Low-mrank approximation

Goal: Approximate $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$ by $\hat{\mathbf{X}}$ such that $\text{mrnk } \hat{\mathbf{X}} \leq (R_1, R_2, R_3)$

Recall that $\hat{\mathbf{X}}$ must admit a Tucker decomposition

$$\hat{\mathbf{X}} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3},$$

where $\mathbf{U}, \mathbf{V}, \mathbf{W}$ are semi-orthogonal w.l.o.g. Hence, a first, natural approach is:

$$\min_{\mathbf{S}, \mathbf{U}, \mathbf{V}, \mathbf{W}} \|\mathbf{X} - (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}\|_{\text{F}}^2.$$

Low-mrank approximation

Goal: Approximate $\mathbf{X} \in \mathbb{R}^{N_1 \times N_2 \times N_3}$ by $\hat{\mathbf{X}}$ such that $\text{mrnk } \hat{\mathbf{X}} \leq (R_1, R_2, R_3)$

Recall that $\hat{\mathbf{X}}$ must admit a Tucker decomposition

$$\hat{\mathbf{X}} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \mathbf{S} \in \mathbb{R}^{R_1 \times R_2 \times R_3},$$

where $\mathbf{U}, \mathbf{V}, \mathbf{W}$ are semi-orthogonal w.l.o.g. Hence, a first, natural approach is:

$$\min_{\mathbf{S}, \mathbf{U}, \mathbf{V}, \mathbf{W}} \|\mathbf{X} - (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}\|_{\text{F}}^2.$$

Since

$$\begin{aligned} \|\mathbf{X} - (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}\|_{\text{F}}^2 &= \|\mathbf{X}\|_{\text{F}}^2 - 2 \langle \mathbf{X}, (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S} \rangle + \|(\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}\|_{\text{F}}^2 \\ &= \|\mathbf{X}\|_{\text{F}}^2 - 2 \left\langle (\mathbf{U}^{\text{T}}, \mathbf{V}^{\text{T}}, \mathbf{W}^{\text{T}}) \cdot \mathbf{X}, \mathbf{S} \right\rangle + \|\mathbf{S}\|_{\text{F}}^2, \end{aligned}$$

it is easy to show that every minimizer must satisfy $\mathbf{S} = (\mathbf{U}^{\text{T}}, \mathbf{V}^{\text{T}}, \mathbf{W}^{\text{T}}) \cdot \mathbf{X}$, and thus the problem is equivalent to

$$\arg \min_{\mathbf{U}, \mathbf{V}, \mathbf{W}} \left\| \mathbf{X} - (\mathbf{U}\mathbf{U}^{\text{T}}, \mathbf{V}\mathbf{V}^{\text{T}}, \mathbf{W}\mathbf{W}^{\text{T}}) \cdot \mathbf{X} \right\|_{\text{F}}^2 = \arg \max_{\mathbf{U}, \mathbf{V}, \mathbf{W}} \overbrace{\left\| (\mathbf{U}^{\text{T}}, \mathbf{V}^{\text{T}}, \mathbf{W}^{\text{T}}) \cdot \mathbf{X} \right\|_{\text{F}}^2}^{\mathbf{S}}$$

Truncated HOSVD (THOSVD)

Recall that the HOSVD of \mathbf{X} is built by computing the SVD of each unfolding:

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \text{where} \quad \begin{cases} \mathbf{X}_{(1)} = \mathbf{U} \Sigma_1 \mathbf{Q}_1^\top \\ \mathbf{X}_{(2)} = \mathbf{V} \Sigma_2 \mathbf{Q}_2^\top \\ \mathbf{X}_{(3)} = \mathbf{W} \Sigma_3 \mathbf{Q}_3^\top \end{cases}$$

Consequence: $\|(\mathbf{S})_{\ell::}\|_F \geq \|(\mathbf{S})_{\ell+1::}\|_F$, and similarly, $\|(\mathbf{S})_{:\ell}\|_F \geq \|(\mathbf{S})_{:\ell+1}\|_F$ and $\|(\mathbf{S})_{::\ell}\|_F \geq \|(\mathbf{S})_{::\ell+1}\|_F$.

Truncated HOSVD (THOSVD)

Recall that the HOSVD of \mathbf{X} is built by computing the SVD of each unfolding:

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}, \quad \text{where} \quad \begin{cases} \mathbf{X}_{(1)} = \mathbf{U} \Sigma_1 \mathbf{Q}_1^\top \\ \mathbf{X}_{(2)} = \mathbf{V} \Sigma_2 \mathbf{Q}_2^\top \\ \mathbf{X}_{(3)} = \mathbf{W} \Sigma_3 \mathbf{Q}_3^\top \end{cases}$$

Consequence: $\|(\mathbf{S})_{\ell::}\|_F \geq \|(\mathbf{S})_{\ell+1::}\|_F$, and similarly, $\|(\mathbf{S})_{:\ell}\|_F \geq \|(\mathbf{S})_{:\ell+1}\|_F$ and $\|(\mathbf{S})_{::\ell}\|_F \geq \|(\mathbf{S})_{::\ell+1}\|_F$.

This suggests truncating¹ \mathbf{U} at R_1 cols, \mathbf{V} at R_2 cols & \mathbf{W} at R_3 cols:

$$\mathbf{X} \approx \hat{\mathbf{X}} = (\mathbf{U}_{:,1:R_1}, \mathbf{V}_{:,1:R_2}, \mathbf{W}_{:,1:R_3}) \cdot \mathbf{S}_{1:R_1,1:R_2,1:R_3}$$

(Special case: $R_1 = R_2 = R_3 = 1$, we get $\mathbf{X} \approx \sigma \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w}$ built from the dominant left singular vectors of each unfolding.)

But is it optimal?

1: De Lathauwer & al., 2000a

Quasi-optimality of truncated HOSVD

Unfortunately, no! Nevertheless, it is said to be **quasi-optimal** in the sense:

Thm (Vannieuwenhoven & al., 2012): Let X^* be a best $\text{mr}(\text{rank}(R_1, \dots, R_d))$ approx. of X , and \hat{X} be obtained by truncating its HOSVD at (R_1, \dots, R_d) . Then,

$$\|X - \hat{X}\|_F^2 \leq d \|X - X^*\|_F^2.$$

Quasi-optimality of truncated HOSVD

Unfortunately, no! Nevertheless, it is said to be **quasi-optimal** in the sense:

Thm (Vannieuwenhoven & al., 2012): Let X^* be a best mrank- (R_1, \dots, R_d) approx. of X , and \hat{X} be obtained by truncating its HOSVD at (R_1, \dots, R_d) . Then,

$$\|X - \hat{X}\|_F^2 \leq d \|X - X^*\|_F^2.$$

Proof ($d = 3$): Define the orthogonal projector $\mathbf{P}_U := \mathbf{U}_{:,1:R_1} \mathbf{U}_{:,1:R_1}^\top$ onto the dominant subspace of dim R_1 of $\mathbf{X}_{(1)}$, and similarly for \mathbf{P}_V and \mathbf{P}_W . Write

$$\begin{aligned} E := X - \hat{X} &= X - (\mathbf{P}_U, \mathbf{P}_V, \mathbf{P}_W) \cdot X = X + (\mathbf{P}_U, \cdot, \cdot) \cdot X - (\mathbf{P}_U, \cdot, \cdot) \cdot X \\ &\quad + (\mathbf{P}_U, \mathbf{P}_V, \cdot) \cdot X - (\mathbf{P}_U, \mathbf{P}_V, \cdot) \cdot X - (\mathbf{P}_U, \mathbf{P}_V, \mathbf{P}_W) \cdot X \\ &= (\mathbf{I} - \mathbf{P}_U, \cdot, \cdot) \cdot X + (\mathbf{P}_U, \mathbf{I} - \mathbf{P}_V, \cdot) \cdot X + (\mathbf{P}_U, \mathbf{P}_V, \mathbf{I} - \mathbf{P}_W) \cdot X \end{aligned}$$

Quasi-optimality of truncated HOSVD

Unfortunately, no! Nevertheless, it is said to be **quasi-optimal** in the sense:

Thm (Vannieuwenhoven & al., 2012): Let \mathbf{X}^* be a best mrank- (R_1, \dots, R_d) approx. of \mathbf{X} , and $\hat{\mathbf{X}}$ be obtained by truncating its HOSVD at (R_1, \dots, R_d) . Then,

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 \leq d \|\mathbf{X} - \mathbf{X}^*\|_F^2.$$

Proof ($d = 3$): Define the orthogonal projector $\mathbf{P}_U := \mathbf{U}_{:,1:R_1} \mathbf{U}_{:,1:R_1}^\top$ onto the dominant subspace of dim R_1 of $\mathbf{X}_{(1)}$, and similarly for \mathbf{P}_V and \mathbf{P}_W . Write

$$\begin{aligned} \mathbf{E} := \mathbf{X} - \hat{\mathbf{X}} &= \mathbf{X} - (\mathbf{P}_U, \mathbf{P}_V, \mathbf{P}_W) \cdot \mathbf{X} = \mathbf{X} + (\mathbf{P}_U, \cdot, \cdot) \cdot \mathbf{X} - (\mathbf{P}_U, \cdot, \cdot) \cdot \mathbf{X} \\ &\quad + (\mathbf{P}_U, \mathbf{P}_V, \cdot) \cdot \mathbf{X} - (\mathbf{P}_U, \mathbf{P}_V, \cdot) \cdot \mathbf{X} - (\mathbf{P}_U, \mathbf{P}_V, \mathbf{P}_W) \cdot \mathbf{X} \\ &= (\mathbf{I} - \mathbf{P}_U, \cdot, \cdot) \cdot \mathbf{X} + (\mathbf{P}_U, \mathbf{I} - \mathbf{P}_V, \cdot) \cdot \mathbf{X} + (\mathbf{P}_U, \mathbf{P}_V, \mathbf{I} - \mathbf{P}_W) \cdot \mathbf{X} \end{aligned}$$

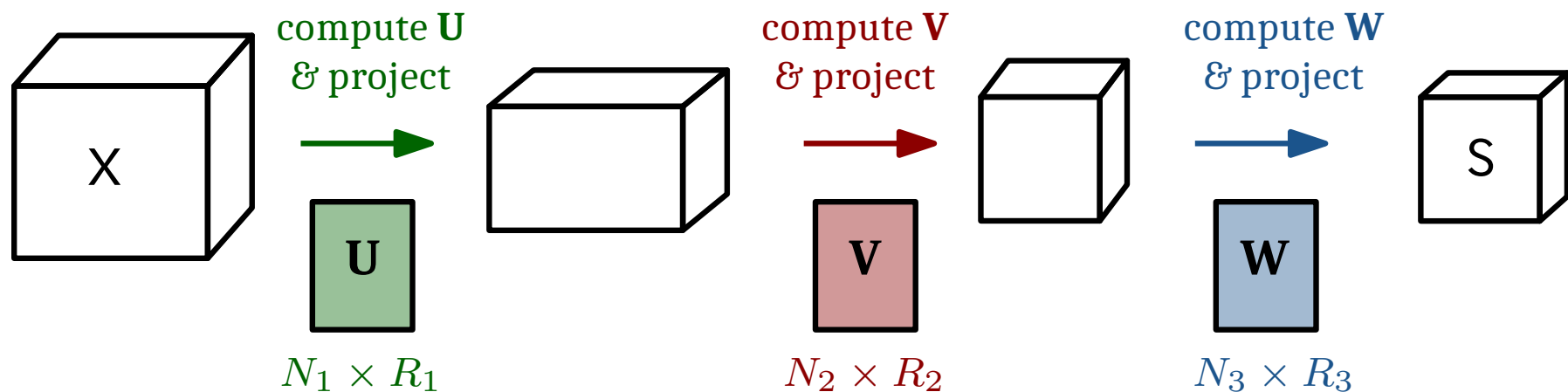
By orthogonality, non-expansiveness and the Eckart–Young thm,

$$\begin{aligned} \|\mathbf{E}\|_F^2 &\leq \|(\mathbf{I} - \mathbf{P}_U, \cdot, \cdot) \cdot \mathbf{X}\|_F^2 + \|(\cdot, \mathbf{I} - \mathbf{P}_V, \cdot) \cdot \mathbf{X}\|_F^2 + \|(\cdot, \cdot, \mathbf{I} - \mathbf{P}_W) \cdot \mathbf{X}\|_F^2 \\ &\leq \sum_{i=1}^3 \left\| \mathbf{X}_{(i)} - \mathbf{X}_{(i)}^* \right\|_F^2 = 3 \|\mathbf{X} - \mathbf{X}^*\|_F^2 \end{aligned}$$



Variants of THOSVD

Variants of THOSVD exist,^{1,2,3} for instance by sequentially performing (optimal) modal projections:



They have lower complexity, often display superior empirical performance & are provably never worse than THOSVD in special cases.

But in general, they're all subject to the same quasi-optimality bound.

1: Vannieuwenhoven & al., 2012, 2: da Silva & al., 2016, 3: Goulart & Comon, 2017

High-order orthogonal iteration (HOOI)

The previous algebraic (non-iterative) soln's can be refined by means of an iterative algorithm.

A popular one is HOOI,¹ which amounts to an alternating opt scheme for

$$\max_{\mathbf{U}, \mathbf{V}, \mathbf{W}} \left\| (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X} \right\|_F^2, \quad \text{subj. to} \quad \begin{cases} \mathbf{U}^\top \mathbf{U} = \mathbf{I}_{R_1} \\ \mathbf{V}^\top \mathbf{V} = \mathbf{I}_{R_2} \\ \mathbf{W}^\top \mathbf{W} = \mathbf{I}_{R_3} \end{cases}$$

1: De Lathauwer & al., 2000b, 2: Xu, 2018

High-order orthogonal iteration (HOOI)

The previous algebraic (non-iterative) soln's can be refined by means of an iterative algorithm.

A popular one is HOOI,¹ which amounts to an alternating opt scheme for

$$\max_{\mathbf{U}, \mathbf{V}, \mathbf{W}} \left\| (\mathbf{U}^\top, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X} \right\|_F^2, \quad \text{subj. to} \quad \begin{cases} \mathbf{U}^\top \mathbf{U} = \mathbf{I}_{R_1} \\ \mathbf{V}^\top \mathbf{V} = \mathbf{I}_{R_2} \\ \mathbf{W}^\top \mathbf{W} = \mathbf{I}_{R_3} \end{cases}$$

For instance, if \mathbf{V} , \mathbf{W} are held fixed, then \mathbf{U} easily follows from a (truncated) SVD. One thus performs at each iteration:

- (i) Recompute \mathbf{U} from the SVD of $((\cdot, \mathbf{V}^\top, \mathbf{W}^\top) \cdot \mathbf{X})_{(1)}$
- (ii) Recompute \mathbf{V} from the SVD of $((\mathbf{U}^\top, \cdot, \mathbf{W}^\top) \cdot \mathbf{X})_{(2)}$
- (iii) Recompute \mathbf{W} from the SVD of $((\mathbf{U}^\top, \mathbf{V}^\top, \cdot) \cdot \mathbf{X})_{(3)}$

Converges to a local sol'n whose unfoldings above have distinct dominating singular values).²

1: De Lathauwer & al., 2000b, 2: Xu, 2018

Best rank- R approximation

Consider now the best rank- R approximation of \mathbf{X} :

$$\inf_{\text{rank } \hat{\mathbf{X}} \leq R} \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F = \inf_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \left\| \mathbf{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]_R \right\|_F$$

Can we replace inf by min? **In general, no: a minimizer might not exist if $R > 1$.**

Best rank- R approximation

Consider now the best rank- R approximation of \mathbf{X} :

$$\inf_{\text{rank } \hat{\mathbf{X}} \leq R} \|\mathbf{X} - \hat{\mathbf{X}}\|_F = \inf_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathbf{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]_R\|_F$$

Can we replace inf by min? **In general, no: a minimizer might not exist if $R > 1$.**

Example: Previously, we have shown that

$$\mathbf{X} = \left(\begin{array}{cc|cc} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{array} \right) = \mathbf{e}_1 \otimes \mathbf{e}_1 \otimes \mathbf{e}_1 + \mathbf{e}_2 \otimes \mathbf{e}_2 \otimes \mathbf{e}_1 + \mathbf{e}_1 \otimes \mathbf{e}_2 \otimes \mathbf{e}_2.$$

has rank 3. More generally, if $\mathbf{A} \in \mathbb{R}^{N_1 \times 2}$, $\mathbf{B} \in \mathbb{R}^{N_2 \times 2}$ and $\mathbf{C} \in \mathbb{R}^{N_3 \times 2}$ all have rank 2, then

$$\mathbf{Y} = (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \mathbf{X} = \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c}_1 + \mathbf{a}_2 \otimes \mathbf{b}_2 \otimes \mathbf{c}_1 + \mathbf{a}_1 \otimes \mathbf{b}_2 \otimes \mathbf{c}_2$$

has rank 3, but does not admit a best rank-2 approximation:

$$\begin{aligned} \hat{\mathbf{Y}}_m &:= m(\mathbf{a}_1 + m^{-1}\mathbf{a}_2) \otimes (\mathbf{b}_2 + m^{-1}\mathbf{b}_1) \otimes (\mathbf{c}_1 + m^{-1}\mathbf{c}_2) - m\mathbf{a}_1 \otimes \mathbf{b}_2 \otimes \mathbf{c}_1 \\ &= \mathbf{Y} + O(1/m) \rightarrow \mathbf{Y} \end{aligned}$$



Ill-posedness of best rank- R approximation

A celebrated paper by de Silva & Lim (2008) discusses this issue in depth.

It argues that there also exist examples of any orders $d \geq 3$ for R in $\{2, \dots, \min_d N_d\}$.

Ill-posedness of best rank- R approximation

A celebrated paper by de Silva & Lim (2008) discusses this issue in depth.

It argues that there also exist examples of any orders $d \geq 3$ for R in $\{2, \dots, \min_d N_d\}$.

Furthermore, by a complete classification of $2 \times 2 \times 2$ tensors by orbit type, they are able to show that in $\mathbb{R}^{N_1 \times N_2 \times N_3}$, no element from the nonempty open set

$$\left\{ (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \left(\begin{array}{cc|cc} 1 & 0 & 0 & -1 \\ 0 & 1 & 1 & 0 \end{array} \right) : \text{rank } \mathbf{A} = \text{rank } \mathbf{B} = \text{rank } \mathbf{C} = 2 \right\}$$

admits a best rank-2 approximation.

Ill-posedness of best rank- R approximation

A celebrated paper by de Silva & Lim (2008) discusses this issue in depth.

It argues that there also exist examples of any orders $d \geq 3$ for R in $\{2, \dots, \min_d N_d\}$.

Furthermore, by a complete classification of $2 \times 2 \times 2$ tensors by orbit type, they are able to show that in $\mathbb{R}^{N_1 \times N_2 \times N_3}$, no element from the nonempty open set

$$\left\{ (\mathbf{A}, \mathbf{B}, \mathbf{C}) \cdot \left(\begin{array}{cc|cc} 1 & 0 & 0 & -1 \\ 0 & 1 & 1 & 0 \end{array} \right) : \text{rank } \mathbf{A} = \text{rank } \mathbf{B} = \text{rank } \mathbf{C} = 2 \right\}$$

admits a best rank-2 approximation.

By contrast, a best rank- R approximation generically exists¹ over \mathbb{C} .

1: Qi & al., 2020

What if we try, anyway?

One might want to dismiss this issue by focusing on a “reasonable rank- R approximation” instead.

What if we try, anyway?

One might want to dismiss this issue by focusing on a “reasonable rank- R approximation” instead.

Yet, if a solution does not exist, then at least some of the terms in the decomposition **must diverge in norm**,¹ so the “approximate solution” is typically useless.

1: de Silva & Lim, 2008, 2: Paatero, 2000

What if we try, anyway?

One might want to dismiss this issue by focusing on a “reasonable rank- R approximation” instead.

Yet, if a solution does not exist, then at least some of the terms in the decomposition **must diverge in norm**,¹ so the “approximate solution” is typically useless.

Also, failing to account for it can bring serious difficulties even when a best approximation with the sought rank exists:

“A well-posed problem in the neighborhood of an ill-posed one is ill-conditioned.”

For instance, a numerical algorithm can traverse regions close to tensors who do not admit a best approximation.²

1: de Silva & Lim, 2008, 2: Paatero, 2000

CPD condition number

Goal: What is the relation between the rank-1 terms of

$$\mathbf{X} = \sum_{r=1}^R \mathbf{X}_r = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$

and those obtained by rank- R approximation $\hat{\mathbf{X}}$ of $\mathbf{Y} \approx \mathbf{X}$ (e.g., $\mathbf{Y} = \mathbf{X} + \mathbf{N}$)?

CPD condition number

Goal: What is the relation between the rank-1 terms of

$$\mathbf{X} = \sum_{r=1}^R \mathbf{X}_r = \sum_{r=1}^R \mathbf{a}_r \otimes \mathbf{b}_r \otimes \mathbf{c}_r$$

and those obtained by rank- R approximation $\hat{\mathbf{X}}$ of $\mathbf{Y} \approx \mathbf{X}$ (e.g., $\mathbf{Y} = \mathbf{X} + \mathbf{N}$)?

Breiding & Vannieuwenhoven (2018) studied the condition number $\kappa(\mathcal{X})$ of the local inverse of

$$\Phi : \mathcal{M} \times \cdots \times \mathcal{M} \rightarrow \mathbb{R}^{N_1 \times N_2 \times N_3}$$

$$\mathcal{X} = (\mathbf{X}_1, \dots, \mathbf{X}_R) \mapsto \sum_{r=1}^R \mathbf{X}_r$$

\mathcal{M} = manifold of rank-1 tensors

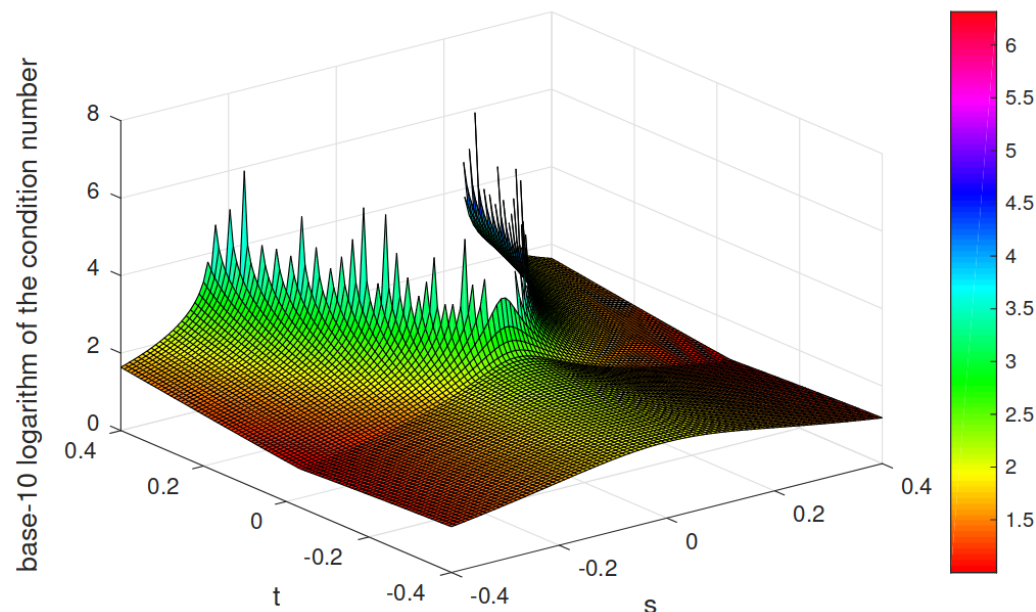
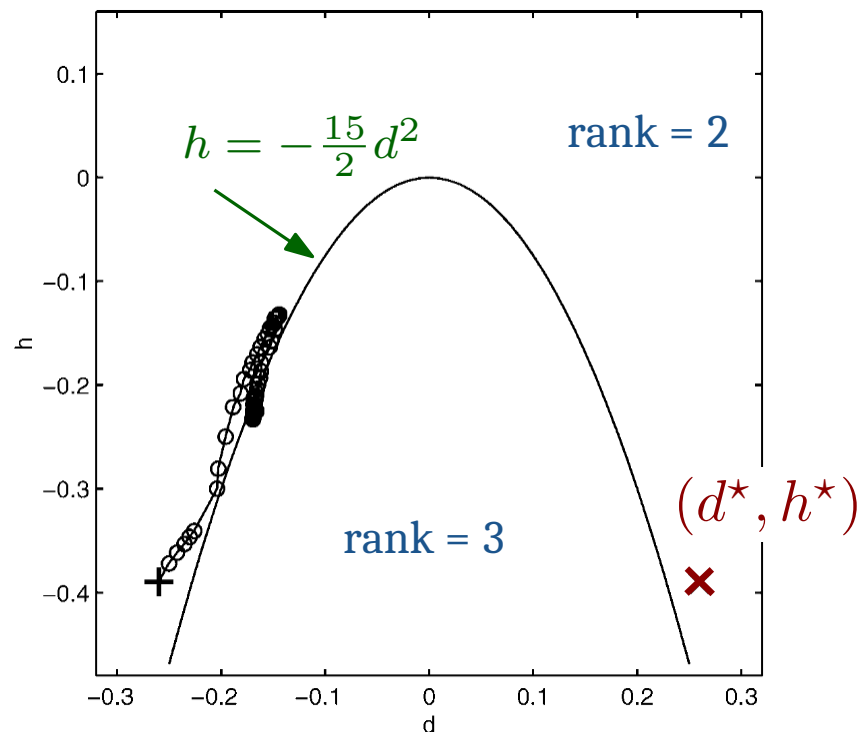
in a neighborhood of \mathcal{X} , denoted $\Phi_{\mathcal{X}}^{-1}$, which allows estimating

$$\underbrace{\left\| \mathcal{X} - \Phi_{\mathcal{X}}^{-1}(\hat{\mathbf{X}}) \right\|}_{\text{error over parameters}} \lesssim \kappa(\mathcal{X}) \underbrace{\left\| \Phi(\mathcal{X}) - \hat{\mathbf{X}} \right\|}_{\text{reconstruction error}},$$

for $\hat{\mathbf{X}}$ a rank- R tensor in the neighborhood of $\mathbf{X} = \Phi(\mathcal{X})$.

Paatero's example

They also showed that $\kappa(\mathcal{X})$ must diverge as $\Phi(\mathcal{X})$ tends to a tensor of rank higher than R .



Corresponding behavior of the condition number (in log scale).

Paatero's example:¹ trajectory of ALS algorithm trying to compute a rank-2 PD of

$$\mathbf{X} = \left(\begin{array}{cc|cc} 0 & 1 & 30 & 0 \\ 1 & d^* & 0 & h^* \end{array} \right).$$

A positive result

In applications, by assumption $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R + \mathbf{W}$, even if \mathbf{X} does not admit a best rank- R approximation.

A positive result

In applications, by assumption $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R + \mathbf{W}$, even if \mathbf{X} does not admit a best rank- R approximation.

For “sufficiently small” \mathbf{W} , the problem is **well posed**:¹

If $R \leq \min \{N_1, N_2\}$, then there exists $\epsilon > 0$ such that \mathbf{X} admits a best rank- R approximation for all \mathbf{W} satisfying $\|\mathbf{W}\|_F \leq \epsilon$.

Roughly, ϵ depends on the conditioning of “partial” CPDs of $\llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$.

1: Evert & De Lathauwer, 2022

A positive result

In applications, by assumption $\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_R + \mathbf{W}$, even if \mathbf{X} does not admit a best rank- R approximation.

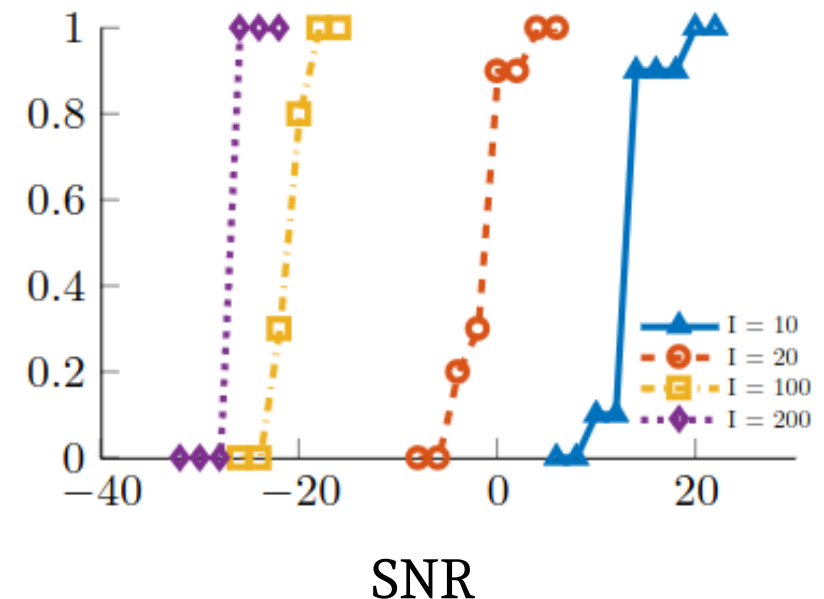
For “sufficiently small” \mathbf{W} , the problem is **well posed**:¹

If $R \leq \min \{N_1, N_2\}$, then there exists $\epsilon > 0$ such that \mathbf{X} admits a best rank- R approximation for all \mathbf{W} satisfying $\|\mathbf{W}\|_F \leq \epsilon$.

Roughly, ϵ depends on the conditioning of “partial” CPDs of $\llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$.

Another, similar result by the same authors allows asserting that a given tensor has a best rank- R approx.

Example:¹ proportion of random tensors $\llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket_4 + \mathbf{W} \in \mathbb{R}^{I \times I \times I}$ guaranteed to have a best rank-4 approx., as a function of the SNR



¹: Evert & De Lathauwer, 2022

Approximate CPD computation

Numerous algorithms exist for computing an approximate CPD, including

- alternating optimization (block coordinate descent)
- algebraic methods
- classical optimization schemes
- stochastic gradient
- distributed schemes
- ...

Several bibliographical pointers are given in the handout.

Best BTD approximation

Recall the BTD model

$$\begin{aligned} \mathbf{X} &\approx \sum_{r=1}^R \mathbf{H}_r \otimes \mathbf{c}_r, \quad \text{rank } \mathbf{H}_r \leq L_r, \\ &= \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{c}_r, \end{aligned}$$

where $\mathbf{A}_r \in \mathbb{R}^{N_1 \times L_r}$, $\mathbf{B}_r \in \mathbb{R}^{N_2 \times L_r}$ and $\mathbf{c}_r \in \mathbb{R}^{N_3}$.

A first, natural formulation would thus be

$$\inf_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{c}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{c}_r \right\|_{\text{F}}^2,$$

with a fixed (given) choice of R and L_1, \dots, L_R .

Does a minimizer always exist? **Not always...**

Sets of BTDs

To study BTD approximations, we need sets analogous to $\{X : \text{rank } X \leq R\}$.

For $L_1 \geq L_2 \geq \dots \geq L_{N_3} \geq 0$, we define:

$$\mathcal{B}_{L_1, \dots, L_{N_3}} \triangleq \left\{ \sum_{r=1}^{N_3} \mathbf{H}_r \otimes \mathbf{w}_r \mid \text{rank } \mathbf{H}_r \leq L_r \text{ and } \mathbf{w}_1, \dots, \mathbf{w}_R \text{ are l. i.} \right\}$$

If $L_{R+1} = \dots = L_{N_3} = 0$, we can simply write $\mathcal{B}_{L_1, \dots, L_R} = \mathcal{B}_{L_1, \dots, L_{N_3}}$

Sets of BTDs

To study BTD approximations, we need sets analogous to $\{X : \text{rank } X \leq R\}$.

For $L_1 \geq L_2 \geq \dots \geq L_{N_3} \geq 0$, we define:

$$\mathcal{B}_{L_1, \dots, L_{N_3}} \triangleq \left\{ \sum_{r=1}^{N_3} \mathbf{H}_r \otimes \mathbf{w}_r \mid \text{rank } \mathbf{H}_r \leq L_r \text{ and } \mathbf{w}_1, \dots, \mathbf{w}_R \text{ are l. i.} \right\}$$

If $L_{R+1} = \dots = L_{N_3} = 0$, we can simply write $\mathcal{B}_{L_1, \dots, L_R} = \mathcal{B}_{L_1, \dots, L_{N_3}}$

Example: If

$$X = \left(\sum_{i=1}^3 \mathbf{u}_i \otimes \mathbf{v}_i \right) \otimes \mathbf{w}_1 + \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{w}_2,$$

then $X \in \mathcal{B}_{3,1} \subset \mathcal{B}_{3,2} \subset \mathcal{B}_{3,2,1} \subset \mathcal{B}_{4,4,2} \quad \square$

Sets of BTDs

To study BTD approximations, we need sets analogous to $\{X : \text{rank } X \leq R\}$.

For $L_1 \geq L_2 \geq \dots \geq L_{N_3} \geq 0$, we define:

$$\mathcal{B}_{L_1, \dots, L_{N_3}} \triangleq \left\{ \sum_{r=1}^{N_3} \mathbf{H}_r \otimes \mathbf{w}_r \mid \text{rank } \mathbf{H}_r \leq L_r \text{ and } \mathbf{w}_1, \dots, \mathbf{w}_R \text{ are l. i.} \right\}$$

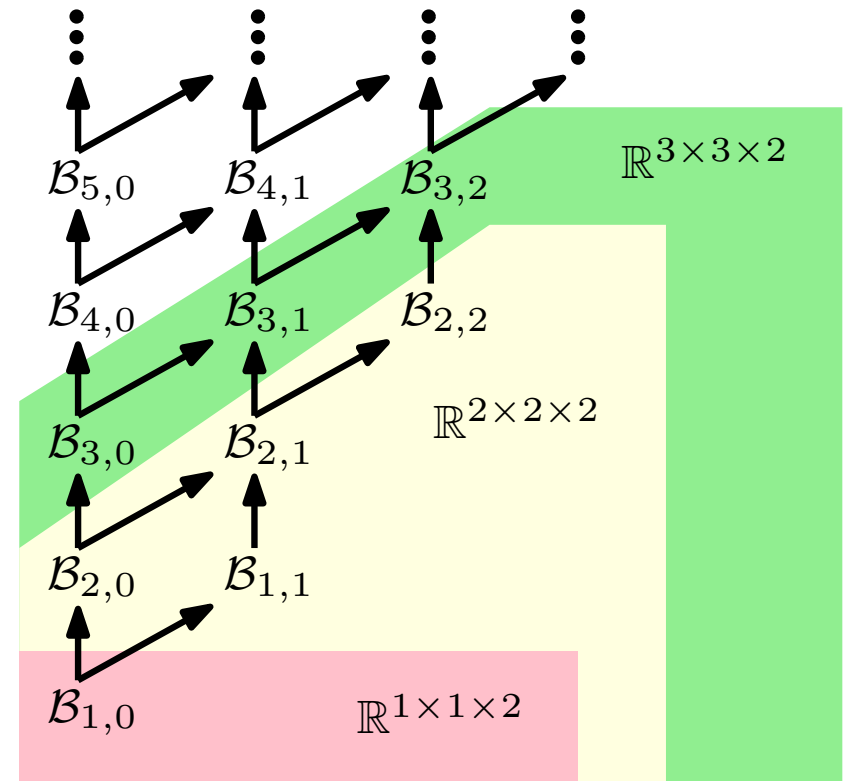
If $L_{R+1} = \dots = L_{N_3} = 0$, we can simply write $\mathcal{B}_{L_1, \dots, L_R} = \mathcal{B}_{L_1, \dots, L_{N_3}}$

Example: If

$$X = \left(\sum_{i=1}^3 \mathbf{u}_i \otimes \mathbf{v}_i \right) \otimes \mathbf{w}_1 + \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{w}_2,$$

then $X \in \mathcal{B}_{3,1} \subset \mathcal{B}_{3,2} \subset \mathcal{B}_{3,2,1} \subset \mathcal{B}_{4,4,2}$ \square

This induces a richer hierarchy than the one induced by the tensor rank.



Hierarchy of BTD sets

Minimal ranks and counter-example

Def: Minimal (BTD) ranks¹ of \mathbf{X}

$$\rho(\mathbf{X}) = (L_1, \dots, L_{N_3}) \quad \text{when} \quad \mathbf{X} \in \mathcal{B}_{S_1, \dots, S_{N_3}} \Leftrightarrow \forall r, S_r \geq L_r.$$

1: Goulart & Comon, 2019.

Minimal ranks and counter-example

Def: Minimal (BTD) ranks¹ of \mathbf{X}

$$\rho(\mathbf{X}) = (L_1, \dots, L_{N_3}) \quad \text{when} \quad \mathbf{X} \in \mathcal{B}_{S_1, \dots, S_{N_3}} \Leftrightarrow \forall r, S_r \geq L_r.$$

Example: $\mathbf{X} \in \mathcal{B}_{4,2} \cap \mathcal{B}_{3,3} \implies \rho(\mathbf{X}) \neq (4, 2) \quad \text{and} \quad \rho(\mathbf{X}) \neq (3, 3).$

□

Now we're ready to generalize the example given for the CPD.

1: Goulart & Comon, 2019.

Minimal ranks and counter-example

Def: Minimal (BTD) ranks¹ of \mathbf{X}

$$\rho(\mathbf{X}) = (L_1, \dots, L_{N_3}) \quad \text{when} \quad \mathbf{X} \in \mathcal{B}_{S_1, \dots, S_{N_3}} \Leftrightarrow \forall r, S_r \geq L_r.$$

Example: $\mathbf{X} \in \mathcal{B}_{4,2} \cap \mathcal{B}_{3,3} \implies \rho(\mathbf{X}) \neq (4, 2) \quad \text{and} \quad \rho(\mathbf{X}) \neq (3, 3).$ □

Now we're ready to generalize the example given for the CPD.

Example: Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{N_1 \times S}$, $\mathbf{C}, \mathbf{D} \in \mathbb{R}^{N_2 \times S}$ and l.i. vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^{N_3}$,

$$\mathbf{X} := (\mathbf{A}\mathbf{C}^\top + \mathbf{B}\mathbf{D}^\top) \otimes \mathbf{v} + (\mathbf{B}\mathbf{C}^\top) \otimes \mathbf{w},$$

$$\mathbf{X}_n := n \left[(\mathbf{B} + n^{-1}\mathbf{A}) (\mathbf{C} + n^{-1}\mathbf{D})^\top \right] \otimes (\mathbf{v} + n^{-1}\mathbf{w}) - n(\mathbf{B}\mathbf{C}^\top) \otimes \mathbf{v} = \mathbf{X} + o(1).$$

$$\searrow \in \mathcal{B}_{S,S}$$

If $\min \{ \text{rank} \begin{pmatrix} \mathbf{A} & \mathbf{B} \end{pmatrix}, \text{rank} \begin{pmatrix} \mathbf{C} & \mathbf{D} \end{pmatrix} \} > R := \frac{3}{2}S$, then $\mathbf{X} \notin \mathcal{B}_{S,S}$, by Sylvester's inequality. Hence:

$$\arg \min_{\hat{\mathbf{X}} \in \mathcal{B}_{S,S}} \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_{\text{F}} = \emptyset.$$

□

Non-existence of best BTD approximation

As in the CPD case, this can happen with positive probability for real tensors:

Thm (G. & Comon): No $2K \times 2K \times 2$ real-valued tensor \mathbf{X} such that $\rho(\mathbf{X}) = (2K, 2K)$ admits a best approximation in $\mathcal{B}_{2K-1, 2K-1}$. These tensors form a non-empty open set.

Non-existence of best BTD approximation

As in the CPD case, this can happen with positive probability for real tensors:

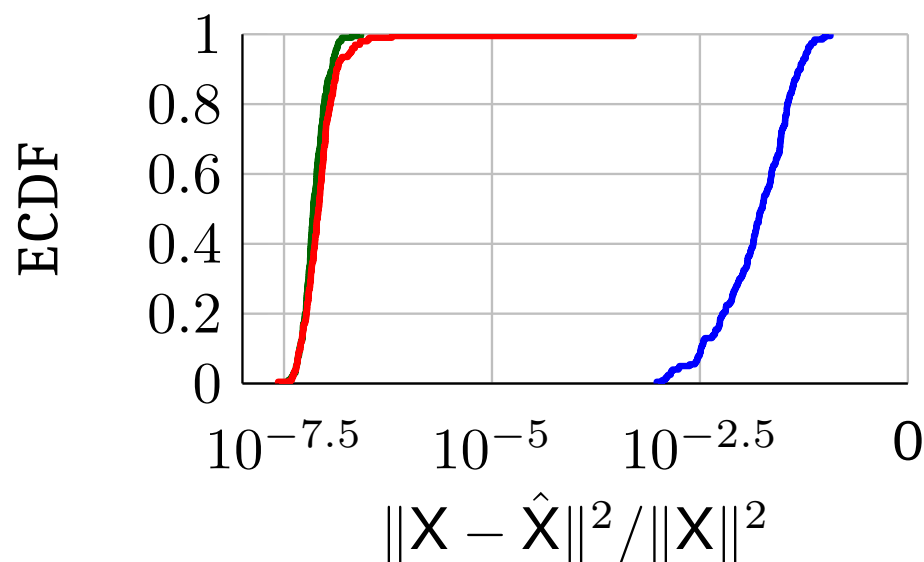
Thm (G. & Comon): No $2K \times 2K \times 2$ real-valued tensor \mathbf{X} such that $\rho(\mathbf{X}) = (2K, 2K)$ admits a best approximation in $\mathcal{B}_{2K-1, 2K-1}$. These tensors form a non-empty open set.

What happens if we try to compute it anyway?

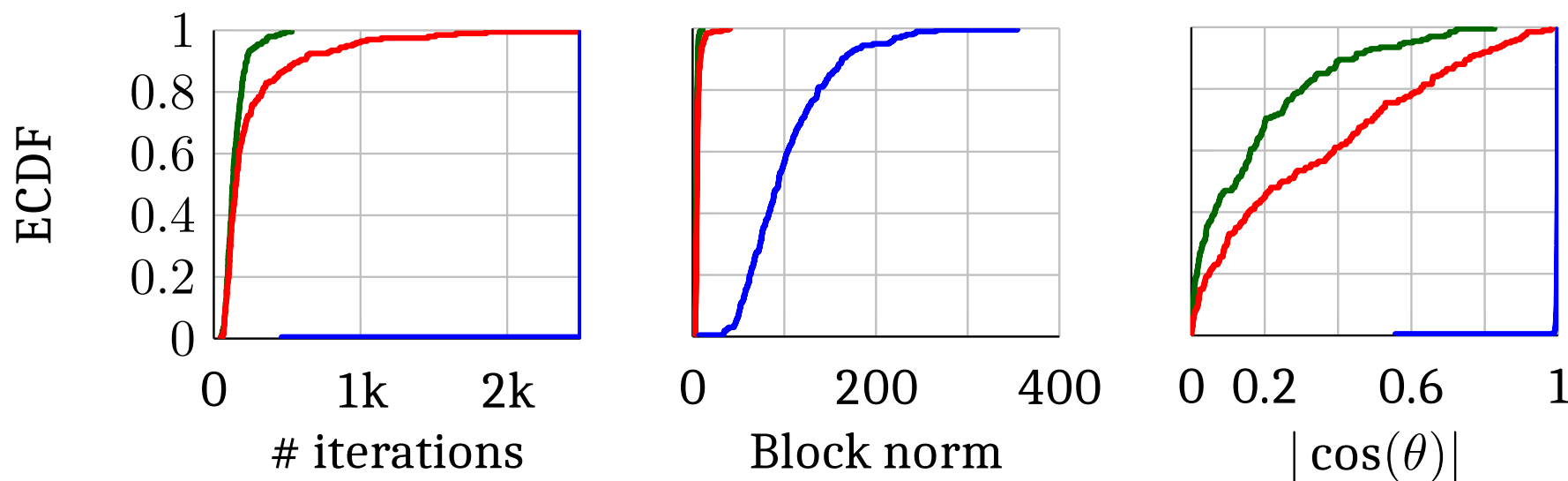
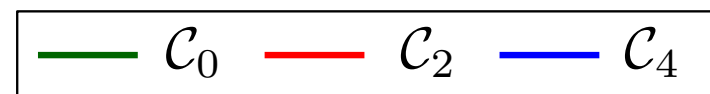
- Optimization algorithm never seems to converge
- Keeps improving error but ultimately gives nearly collinear blocks
- Norms of these blocks “blow up,” but overall error stays bounded

Numerical example

Approximation in $\mathcal{B}_{3,3}$ of $4 \times 4 \times 2$ tensor with opt. algorithm (best of 50 init.):



Classes of tensors, as per # of eigenvalues $\in \mathbb{C}$ of associated pencil:



Computation of approximate BTD

Several algorithms exist (see handout), including:

- alternating optimization (BCD)
- algebraic methods
- standard optimization schemes

Computation of approximate BTD

Several algorithms exist (see handout), including:

- alternating optimization (BCD)
- algebraic methods
- standard optimization schemes

Some of them are based on the standard least-squares formulation

$$\inf_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{C}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{C}_r \right\|_F^2.$$

Problems with this approach:

- the structure (L_1, \dots, L_R) must be fixed a priori
- a solution might not exist
- traversing regions of ill-conditioned BTDs (& thus slow progress)
- poor local minima due to **block rank inversion**

Joint estimation of BTD parameters & ranks

Alternative: automatic rank selection by regularization

$$\min_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{C}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{C}_r \right\|_{\text{F}}^2 + \lambda \underbrace{\sum_{r=1}^R \left(\sum_{m=1}^{L_m} (\|\mathbf{A}_r\|_{:,m} + \|\mathbf{B}_r\|_{:,m}) + \|\mathbf{C}_r\| \right)}_{= \|\mathbf{A}\|_{2,1} + \|\mathbf{B}\|_{2,1} + \|\mathbf{C}\|_{2,1}}$$

Solutions always exist (coercive objective).

Joint estimation of BTD parameters & ranks

Alternative: automatic rank selection by regularization

$$\min_{\{\mathbf{A}_r, \mathbf{B}_r, \mathbf{C}_r\}_{r=1}^R} \left\| \mathbf{X} - \sum_{r=1}^R \left(\mathbf{A}_r \mathbf{B}_r^\top \right) \otimes \mathbf{C}_r \right\|_F^2 + \lambda \underbrace{\sum_{r=1}^R \left(\sum_{m=1}^{L_m} (\|\mathbf{A}_r\|_{:,m} + \|\mathbf{B}_r\|_{:,m}) + \|\mathbf{C}_r\| \right)}_{= \|\mathbf{A}\|_{2,1} + \|\mathbf{B}\|_{2,1} + \|\mathbf{C}\|_{2,1}}$$

Solutions always exist (coercive objective).

Alternating group lasso (AGL) algorithm¹:

- block-coordinate descent algorithm
- each subproblem is a (convex) group lasso problem
- by adding a proximal term, all limit points are stationary points²

Subproblem in \mathbf{A} (those in \mathbf{B} and \mathbf{C} are similar):

$$\min_{\mathbf{A}} \|\text{vec}(\mathbf{X}) - \mathbf{W}_{\mathbf{B}, \mathbf{C}} \text{vec}(\mathbf{A})\|^2 + \lambda \|\mathbf{A}\|_{2,1} + \tau \|\mathbf{A} - \mathbf{A}_0\|_F^2$$

1: Goulart & al., 2020, 2: Razaviyayn & al., 2013

Tensor PCA & asymptotic MLE performance



Tensor PCA & large-dimensional regime

Large body of recent work on the **tensor PCA** problem¹

$$\max_{\|\mathbf{u}\|=1} \sum_{ijk} x_{ijk} u_i u_j u_k = \max_{\|\mathbf{u}\|=1} \lambda \langle \mathbf{a}, \mathbf{u} \rangle^3 + \sum_{ijk} w_{ijk} u_i u_j u_k$$

since the introduction of the spiked (symmetric) rank-1 model (same as before)

$$\mathbf{X} = \lambda \mathbf{a}^{\otimes 3} + \mathbf{W}.$$

Special attention is paid to the **large-dimensional regime**, as $N \rightarrow \infty$.

1: Montanari & Richard, 2014

Tensor PCA & large-dimensional regime

Large body of recent work on the **tensor PCA** problem¹

$$\max_{\|\mathbf{u}\|=1} \sum_{ijk} x_{ijk} u_i u_j u_k = \max_{\|\mathbf{u}\|=1} \lambda \langle \mathbf{a}, \mathbf{u} \rangle^3 + \sum_{ijk} w_{ijk} u_i u_j u_k$$

since the introduction of the spiked (symmetric) rank-1 model (same as before)

$$\mathbf{X} = \lambda \mathbf{a}^{\otimes 3} + \mathbf{W}.$$

Special attention is paid to the **large-dimensional regime**, as $N \rightarrow \infty$.

Connections to many fields & exciting hot topics (see refs. on handout):

- study of disordered systems, spin glasses & statistical physics
- random optimization landscape
- high-dimensional probability & statistics
- random matrix theory

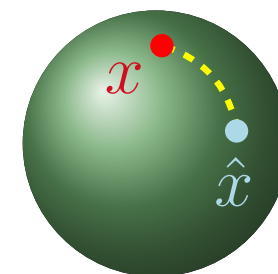
1: Montanari & Richard, 2014

Asymptotic performance limits?

Given any estimator $\hat{\mathbf{a}} : \mathcal{S}^3(N) \rightarrow \mathbb{S}^{N-1}$, a natural performance measure is:

Def: alignment (or overlap)

$$\alpha_{3,N}(\lambda) := \langle \mathbf{a}, \hat{\mathbf{a}}(\mathbf{X}) \rangle \in [-1, 1]$$



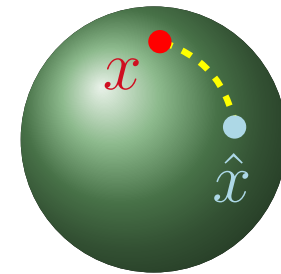
Example: If $\hat{\mathbf{a}} \sim \mathcal{U}(\mathbb{S}^{N-1})$, then asymptotically $\mathbf{a} \perp \hat{\mathbf{a}}$ almost surely.

Asymptotic performance limits?

Given any estimator $\hat{a} : \mathcal{S}^3(N) \rightarrow \mathbb{S}^{N-1}$, a natural performance measure is:

Def: alignment (or overlap)

$$\alpha_{3,N}(\lambda) := \langle \mathbf{a}, \hat{\mathbf{a}}(\mathbf{X}) \rangle \in [-1, 1]$$



Example: If $\hat{a} \sim \mathcal{U}(\mathbb{S}^{N-1})$, then asymptotically $\mathbf{a} \perp \hat{\mathbf{a}}$ almost surely.

Central questions:

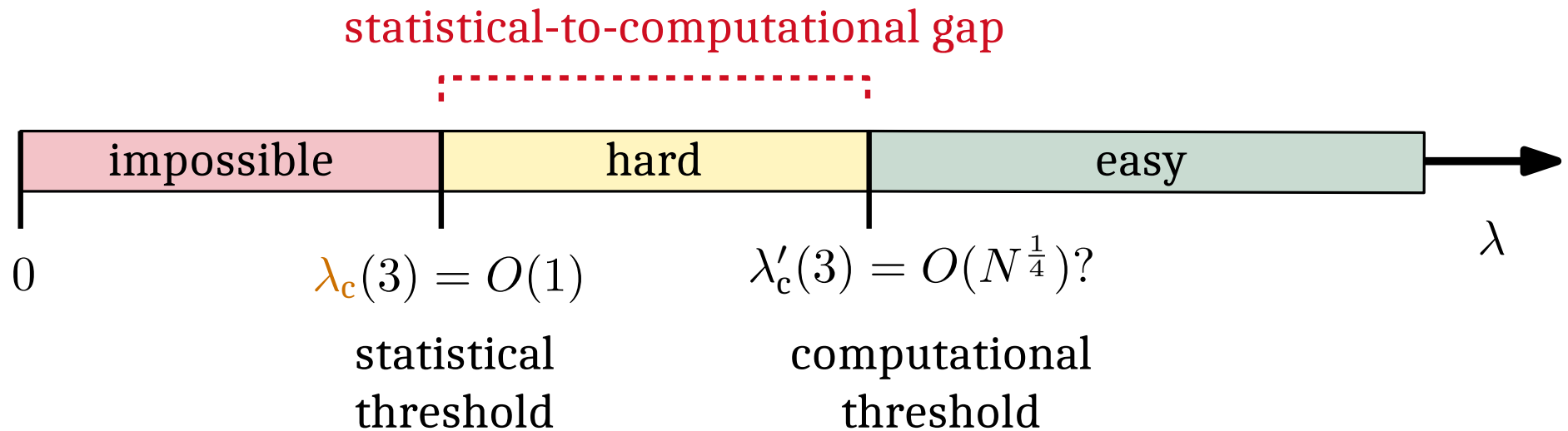
1. **Weak recovery:** for which range of λ is there a \hat{a} such that

$$\limsup_{N \rightarrow \infty} \mathbb{E} \{ \alpha_{3,N}(\lambda) \} > 0 ?$$

2. **Best asymptotic alignment:** what is the largest attainable value of $\limsup_{N \rightarrow \infty} \mathbb{E} \{ \alpha_{3,N}(\lambda) \}$ for each λ ?

Answers and conjectured gap

1. Regimes of weak recovery:



2. Maximum likelihood estimation (MLE) attains the information-theoretic bound on the alignment for all λ .

MLE performance

Settled by Jagannath–Lopatto–Miolane (2020), thanks to [spin glass theory](#):

$$\mu_{d,N}^*(\lambda) = \max_{\|\mathbf{u}\|=1} \left\{ \lambda \langle \mathbf{a}, \mathbf{u} \rangle^d + \mathbf{W} \cdot \mathbf{u}^d \right\} \xrightarrow{\text{a.s.}} \text{GS}_d + \int_0^\lambda q_d^*(t)^{d/2} dt$$

$$|\alpha_{d,N}(\lambda)| = |\langle \mathbf{a}, \hat{\mathbf{a}} \rangle| \xrightarrow{\text{a.s.}} \sqrt{q_d^*(\lambda)}$$

Explicit expressions exist for $d = 3, 4, 5$.

MLE performance

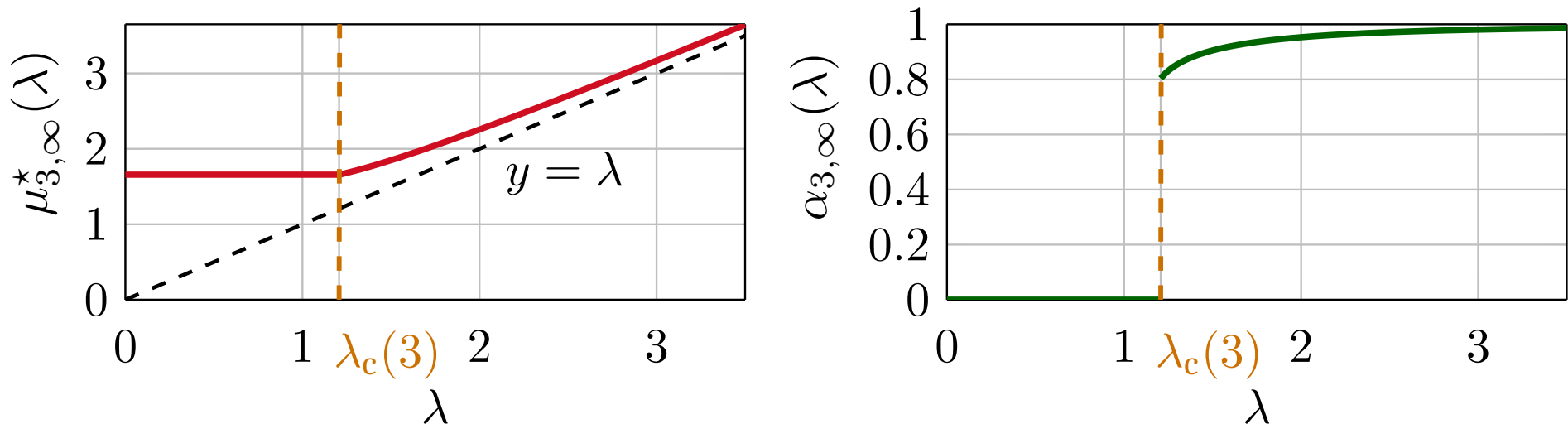
Settled by Jagannath–Lopatto–Miolane (2020), thanks to [spin glass theory](#):

$$\mu_{d,N}^*(\lambda) = \max_{\|\mathbf{u}\|=1} \left\{ \lambda \langle \mathbf{a}, \mathbf{u} \rangle^d + \mathbf{W} \cdot \mathbf{u}^d \right\} \xrightarrow{\text{a.s.}} \text{GS}_d + \int_0^\lambda q_d^*(t)^{d/2} dt$$

$$|\alpha_{d,N}(\lambda)| = |\langle \mathbf{a}, \hat{\mathbf{a}} \rangle| \xrightarrow{\text{a.s.}} \sqrt{q_d^*(\lambda)}$$

Explicit expressions exist for $d = 3, 4, 5$.

For all d , these quantities undergo a phase transition at a threshold $\lambda_c(d) = O(1)$.



Furthermore, the MLE attains the bound $\limsup_N \mathbb{E} \{ |\langle \mathbf{a}, \hat{\mathbf{a}} \rangle| \} \leq \sqrt{q_d^*(\lambda)}$

Extension to other spiked models?

However, it is not obvious how to use these tools to handle other, more general, models.

This motivated our recent contribution¹ where we carry out a similar analysis using tools from [random matrix theory](#).

Using that approach, the estimation of other models has been addressed, notably in the asymmetric case:²

$$X = \lambda \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} + W$$

1: Goulart & al., 2022, 2: Seddik & al., 2024, 3: Lebeau & al., 2024b

Extension to other spiked models?

However, it is not obvious how to use these tools to handle other, more general, models.

This motivated our recent contribution¹ where we carry out a similar analysis using tools from [random matrix theory](#).

Using that approach, the estimation of other models has been addressed, notably in the asymmetric case:²

$$\mathbf{X} = \lambda \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} + \mathbf{W}$$

and for a nested matrix-tensor model³ which applies to a simplified multi-view clustering model:

$$\mathbf{X} = (\boldsymbol{\mu} \mathbf{b}^\top + \mathbf{Z}) \otimes \mathbf{c} + \mathbf{W},$$

where $\boldsymbol{\mu} \in \mathbb{R}^p$ models cluster means, $\mathbf{b} \in \{-1, 1\}^N$ models cluster assignments, \mathbf{Z} is a Gaussian matrix modeling data dispersion, \mathbf{W} models measurement noise and $\mathbf{c} \in \mathbb{R}^M$ models varying SNR conditions.

1: Goulart & al., 2022, 2: Seddik & al., 2024, 3: Lebeau & al., 2024b

Tensor and matrix eigenpairs

Another characterization of tensor eigenpairs (assuming $\|\mathbf{u}\| = 1$):

$$(\mu, \mathbf{u}) \text{ eigenpair of } \mathbf{Y} \quad \Leftrightarrow \quad (\mu, \mathbf{u}) \text{ eigenpair of } \mathbf{Y} \cdot \mathbf{u}^{d-2}$$

Proof: $\mu \mathbf{u} = \mathbf{Y} \cdot \mathbf{u}^{d-1} = (\mathbf{Y} \cdot \mathbf{u}^{d-2}) \mathbf{u}$



Tensor and matrix eigenpairs

Another characterization of tensor eigenpairs (assuming $\|\mathbf{u}\| = 1$):

$$(\mu, \mathbf{u}) \text{ eigenpair of } \mathbf{Y} \quad \Leftrightarrow \quad (\mu, \mathbf{u}) \text{ eigenpair of } \mathbf{Y} \cdot \mathbf{u}^{d-2}$$

Proof: $\mu \mathbf{u} = \mathbf{Y} \cdot \mathbf{u}^{d-1} = (\mathbf{Y} \cdot \mathbf{u}^{d-2}) \mathbf{u}$



Observed performance of algorithms

The performance of power iteration evidently depends on its initialization.

In tensor PCA, randomly initialized power iteration has^{1,2} a (conjectured) **algorithmic threshold** of $\lambda_{\text{alg}} = O(N^{\frac{d-1}{2}})$. This is the SNR required to “beat entropy.”

With a spectral initialization (unfolding SVD), its threshold is³ $\lambda'_{\text{alg}} = O(N^{\frac{d-2}{4}})$.

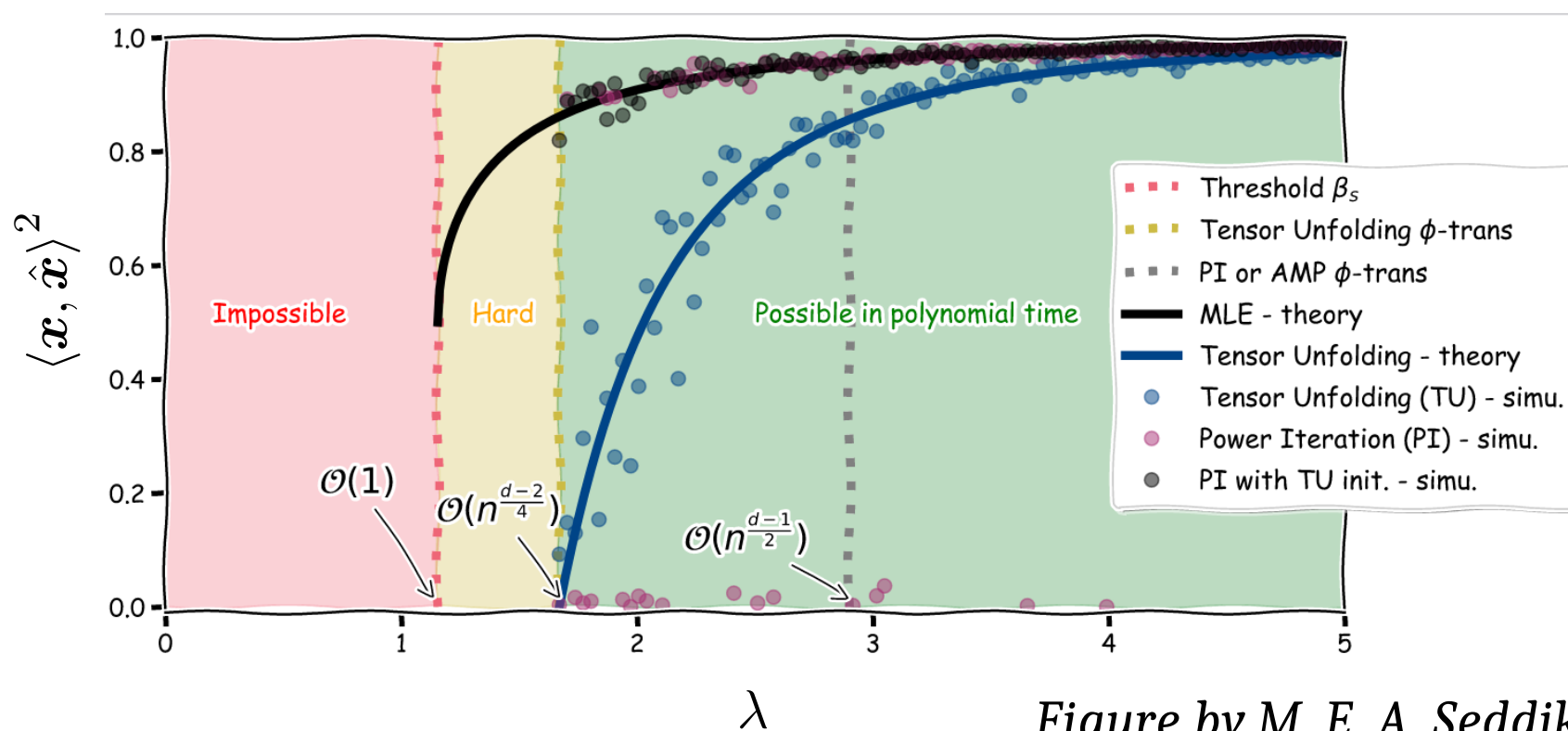
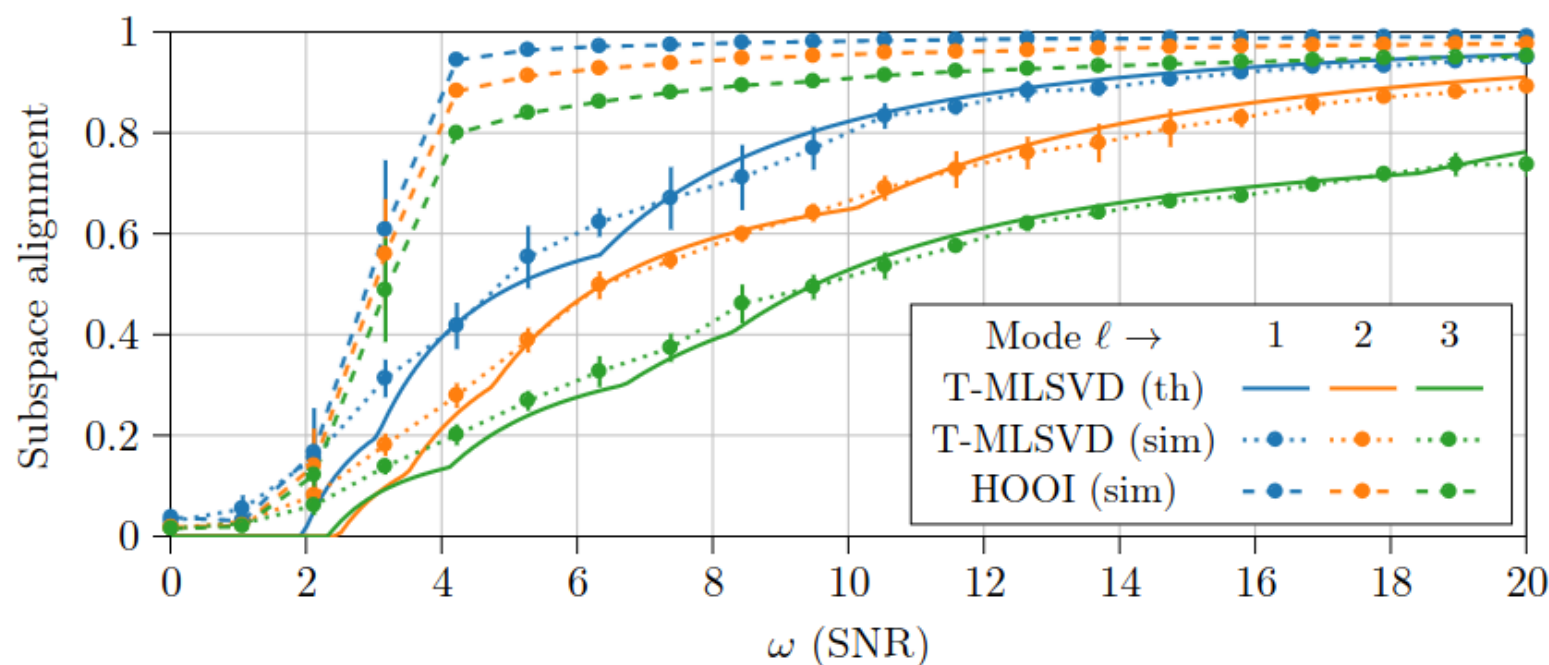


Figure by M. E. A. Seddik.

Performance of THOSVD & HOOI

The large- N performance of THOSVD was analyzed under Gaussian noise by Lebeau & al. (2024a), unveiling a phase transition w.r.t. the SNR at $O(N^{\frac{d-2}{4}})$.

Similar theoretical results for the MLE (best low-mrank approximation) do not yet exist, only empirical ones, eg:



Performances of THOSVD and HOOI¹ for $(N_1, N_2, N_3) = (100, 200, 300)$ and $(R_1, R_2, R_3) = (3, 4, 5)$.

¹: Lebeau & al., 2024a

See you next fall !



Thematic Trimester

Beyond classical regimes in statistical inference and machine learning

 September to November 2024

 Institut de Mathématiques de
Toulouse - France

1 Colloquium
2 Thematic Schools &
1 Workshop,

- **Opening Colloquium**
September 11th, 2024
- **Thematic School: Optimization & algorithms for
high-dimensional machine learning and inference**
October 7th to 11th, 2024
- **Thematic School: Models & methods for high-
dimensional machine learning and inference**
October 14th to 18th, 2024
- **Workshop**
November 4th to 8th, 2024

Organizers: Henrique Goulart (IRIT/Toulouse INP), Vanessa Kientz (CEA List), Vincent Lahoche (CEA List), Xiaoyi Mai (IMT/UT2J), Mohamed Tamaazousti (CEA List)

<https://indico.math.cnrs.fr/category/682>

cimi
Centre International de Mathématiques et d'Informatique
TOULOUSE



Supplementary slides

Tensors: (sketch of a) formal definition

The idea is to impose the multilinear structure by means of a quotient space.

Take, for instance, two vector spaces $\mathcal{U} = \mathbb{R}^{N_1}$ and $\mathcal{V} = \mathbb{R}^{N_2}$. For any pair $(\mathbf{u}, \mathbf{v}) \in \mathcal{U} \times \mathcal{V}$, we would like that

$$\mathbf{u} \otimes \mathbf{v} = (\alpha \mathbf{u}) \otimes (\alpha^{-1} \mathbf{v}), \quad \forall \alpha \neq 0.$$

Hence, $\mathbf{u} \otimes \mathbf{v}$ can be seen as an equivalence class defined on the **free vector space** $F(\mathcal{U} \times \mathcal{V})$ containing all elements $f_{(\alpha \mathbf{u}, \alpha^{-1} \mathbf{v})} \in F(\mathcal{U} \times \mathcal{V})$.

This leads to the definition: $\mathcal{U} \otimes \mathcal{V} := F(\mathcal{U} \times \mathcal{V})/\mathcal{N}$, where \mathcal{N} is the “null subspace”

$$\mathcal{N} := \text{span} \left\{ \sum_{i_1} \sum_{i_2} \alpha_{i_1} \beta_{i_2} f(\mathbf{u}_{i_1}, \mathbf{v}_{i_2}) - f\left(\sum_{i_1} \alpha_{i_1} \mathbf{u}_{i_1}, \sum_{i_2} \beta_{i_2} \mathbf{v}_{i_2}\right) : \mathbf{u}_{i_1} \in \mathcal{U}, \mathbf{v}_{i_2} \in \mathcal{V} \right\}.$$

Example #5: High-dim density approximation

Recent works have relied on low-rank tensor approximation for representing (& learning) densities in high-dim spaces, such as:¹

$$q(\mathbf{x}) = \sum_{n_1=1}^N \cdots \sum_{n_d=1}^N a_{n_1, \dots, n_d} f_{n_1}(x_1) \cdots f_{n_d}(x_d) = \langle \mathbf{A}, \mathbf{F}(\mathbf{x}) \rangle,$$

with $\mathbf{F}(\mathbf{x}) = \mathbf{f}(x_1) \otimes \cdots \otimes \mathbf{f}(x_d)$ and $\mathbf{f}(x_i) = (f_1(x_i) \cdots f_N(x_i))^T$ a vector of chosen basis functions evaluated at x_i .

Example #5: High-dim density approximation

Recent works have relied on low-rank tensor approximation for representing (& learning) densities in high-dim spaces, such as:¹

$$\begin{array}{c} \in \mathbb{R}^d \\ \vdots \\ q(\mathbf{x}) \end{array} = \sum_{n_1=1}^N \cdots \sum_{n_d=1}^N a_{n_1, \dots, n_d} f_{n_1}(x_1) \cdots f_{n_d}(x_d) = \langle \mathbf{A}, \mathbf{F}(\mathbf{x}) \rangle,$$

with $\mathbf{F}(\mathbf{x}) = \mathbf{f}(x_1) \otimes \cdots \otimes \mathbf{f}(x_d)$ and $\mathbf{f}(x_i) = (f_1(x_i) \ \cdots \ f_N(x_i))^T$ a vector of chosen basis functions evaluated at x_i .

The curse of dimensionality is broken by imposing a **low-rank structure** on \mathbf{A} , namely, a **tensor-train structure**:

$$\begin{array}{ccccccc} a_{n_1, \dots, n_d} & = & \mathbf{g}^{(1)}(n_1)^T & \mathbf{G}^{(2)}(n_2) & \cdots & \mathbf{G}^{(d-1)}(n_{d-1}) & \mathbf{g}^{(d)}(n_d). \\ & & \vdots & \vdots & & \vdots & \vdots \\ & & R_1 & R_1 \times R_2 & R_{d-2} \times R_{d-1} & R_{d-1} \end{array}$$

The coordinate tensor \mathbf{A} (of dim N^d) is parameterized by the above vectors & matrices, whose sizes control the model complexity.

1: Novikov & al., 2021

Pointwise evaluation, marginalizing & sampling

The separability in x_1, \dots, x_d is key to complexity reduction:

$$q(\mathbf{x}) = \left(\sum_{n_1} f_{n_1}(\mathbf{x}_1) \mathbf{g}^{(1)}(n_1) \right)^\top \left(\sum_{n_2} f_{n_2}(\mathbf{x}_2) \mathbf{g}^{(2)}(n_2) \right) \dots \left(\sum_{n_d} f_{n_d}(\mathbf{x}_d) \mathbf{g}^{(d)}(n_d) \right)$$

which takes $O(dNR^2)$ operations if $R_i = R$ for all i , instead of N^d !

Pointwise evaluation, marginalizing & sampling

The separability in x_1, \dots, x_d is key to complexity reduction:

$$q(\mathbf{x}) = \left(\sum_{n_1} \mathbf{f}_{n_1}(\mathbf{x}_1) \mathbf{g}^{(1)}(n_1) \right)^\top \left(\sum_{n_2} \mathbf{f}_{n_2}(\mathbf{x}_2) \mathbf{G}^{(2)}(n_2) \right) \dots \left(\sum_{n_d} \mathbf{f}_{n_d}(\mathbf{x}_d) \mathbf{g}^{(d)}(n_d) \right)$$

which takes $O(dNR^2)$ operations if $R_i = R$ for all i , instead of N^d !

Similarly, marginals can be computed using

$$\begin{aligned} q(x_1, \dots, x_{k-1}) &= \int q(\mathbf{x}) \, dx_k \dots dx_d \\ &= \left\langle \mathbf{A}, \mathbf{f}(x_1) \otimes \dots \otimes \mathbf{f}(x_{k-1}) \otimes \left(\int \mathbf{f}(x_k) \, dx_k \right) \otimes \dots \otimes \left(\int \mathbf{f}(x_d) \, dx_d \right) \right\rangle \end{aligned}$$

Pointwise evaluation, marginalizing & sampling

The separability in x_1, \dots, x_d is key to complexity reduction:

$$q(\mathbf{x}) = \left(\sum_{n_1} \mathbf{f}_{n_1}(\mathbf{x}_1) \mathbf{g}^{(1)}(n_1) \right)^\top \left(\sum_{n_2} \mathbf{f}_{n_2}(\mathbf{x}_2) \mathbf{G}^{(2)}(n_2) \right) \dots \left(\sum_{n_d} \mathbf{f}_{n_d}(\mathbf{x}_d) \mathbf{g}^{(d)}(n_d) \right)$$

which takes $O(dNR^2)$ operations if $R_i = R$ for all i , instead of N^d !

Similarly, marginals can be computed using

$$\begin{aligned} q(x_1, \dots, x_{k-1}) &= \int q(\mathbf{x}) \, dx_k \dots dx_d \\ &= \left\langle \mathbf{A}, \mathbf{f}(x_1) \otimes \dots \otimes \mathbf{f}(x_{k-1}) \otimes \left(\int \mathbf{f}(x_k) \, dx_k \right) \otimes \dots \otimes \left(\int \mathbf{f}(x_d) \, dx_d \right) \right\rangle \end{aligned}$$

(Conditional) CDFs can also be computed with 1D integration and summations/matrix-vector products, allowing in particular **efficient sampling** from $q(\mathbf{x})$.

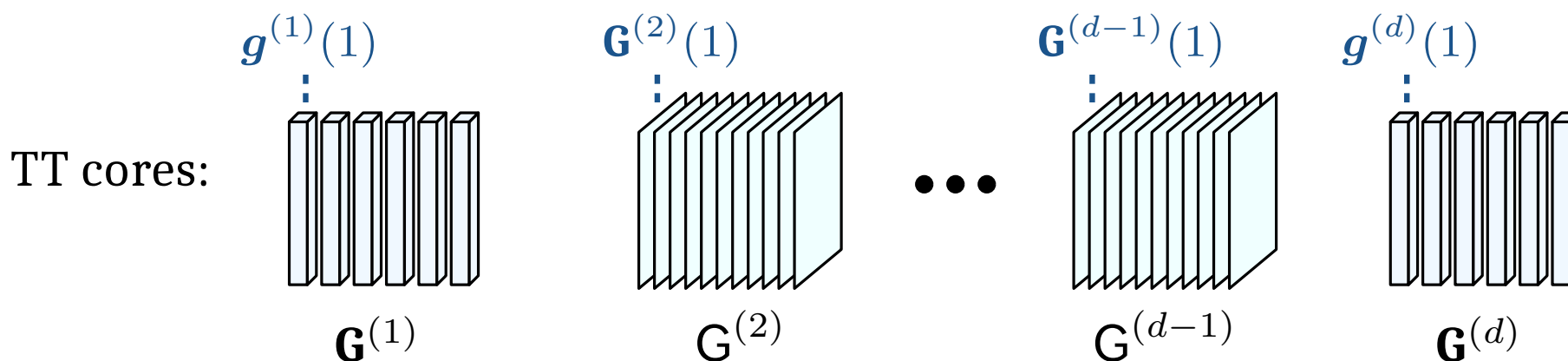
The tensor train (TT) model

Def: Tensor train decomposition¹ of $\mathbf{X} \in \mathbb{R}^{N_1 \times \dots \times N_d}$:

$$x_{n_1, \dots, n_d} = \underset{\substack{\vdots \\ R_1}}{\mathbf{g}^{(1)}(n_1)}^\top \underset{\substack{\vdots \\ R_1 \times R_2}}{\mathbf{G}^{(2)}(n_2)} \dots \underset{\substack{\vdots \\ R_{d-2} \times R_{d-1}}}{\mathbf{G}^{(d-1)}(n_{d-1})} \underset{\substack{\vdots \\ R_{d-1}}}{\mathbf{g}^{(d)}(n_d)}$$

parameterized by **core tensors** $\mathbf{G}^{(1)} \in \mathbb{R}^{R_1 \times N_1}$, $\mathbf{G}^{(i)} \in \mathbb{R}^{R_{i-1} \times N_i \times R_i}$, $i = 2, \dots, d-1$, and $\mathbf{G}^{(d)} \in \mathbb{R}^{R_{d-1} \times N_d}$.

The **TT-rank** of \mathbf{X} is (R_1, \dots, R_{d-1}) .

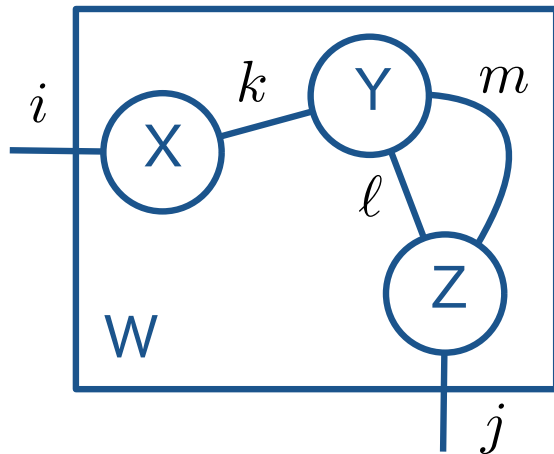


Tensor contractions

Contraction among tensors = summation over some dimensions.

Useful pictorial notation: **tensor networks**¹ displaying tensors as nodes and edges as their indices. A connection then means a contraction.

Ex:



$$w_{ij} = \sum_k \sum_\ell \sum_m x_{ik} y_{k\ell m} z_{\ell m j}$$

Def: Scalar product

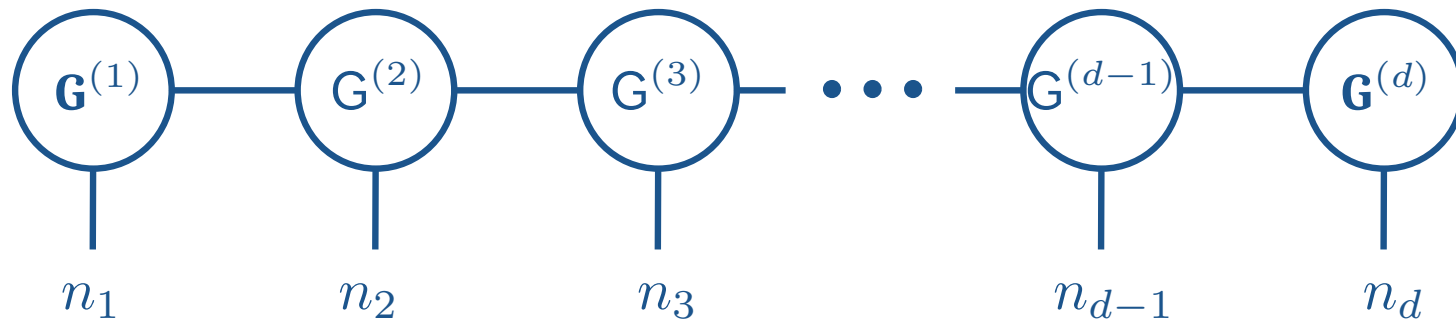
$$\langle X, Y \rangle = \sum_{ijk} x_{ijk} y_{ijk}$$



In particular:
$$\begin{cases} \langle a \otimes b \otimes c, Y \rangle = (a, b, c) \cdot Y \\ \langle a \otimes b \otimes c, u \otimes v \otimes w \rangle = \langle a, u \rangle \langle b, v \rangle \langle c, w \rangle \end{cases}$$

Properties of TT

The TT belongs to a larger class of [hierarchical Tucker models](#)¹, and can be seen as a tensor network of the form

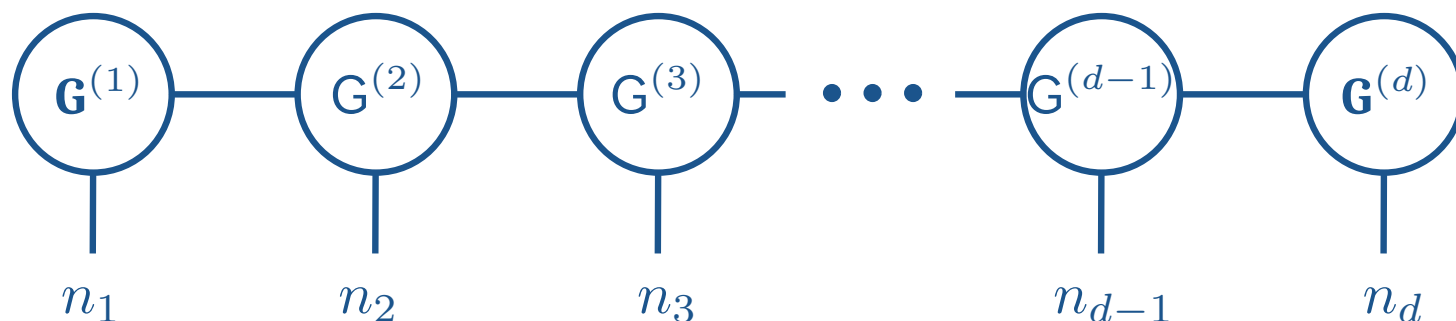


It is also known as the [matrix-product state \(MPS\)](#) model in the physics community.²

1: Hackbusch, 2012, 2: Vidal, 2003

Properties of TT

The TT belongs to a larger class of **hierarchical Tucker models**¹, and can be seen as a tensor network of the form



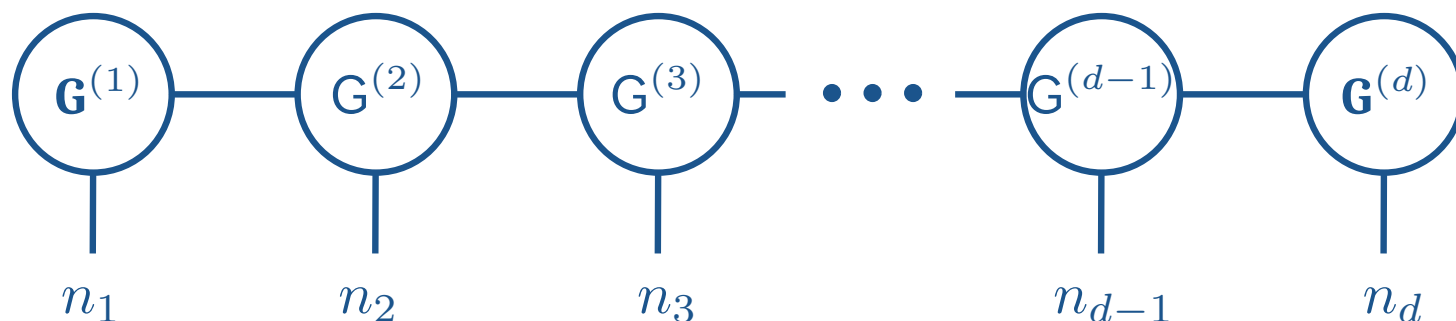
It is also known as the **matrix-product state (MPS)** model in the physics community.²

Motivation: **complexity reduction**. In particular, if $N_i = N$ and $R_i = R$, reduces storage cost from N^d to $O(dN^2R)$.

1: Hackbusch, 2012, 2: Vidal, 2003

Properties of TT

The TT belongs to a larger class of **hierarchical Tucker models**¹, and can be seen as a tensor network of the form



It is also known as the **matrix-product state (MPS)** model in the physics community.²

Motivation: **complexity reduction**. In particular, if $N_i = N$ and $R_i = R$, reduces storage cost from N^d to $O(dN^2R)$.

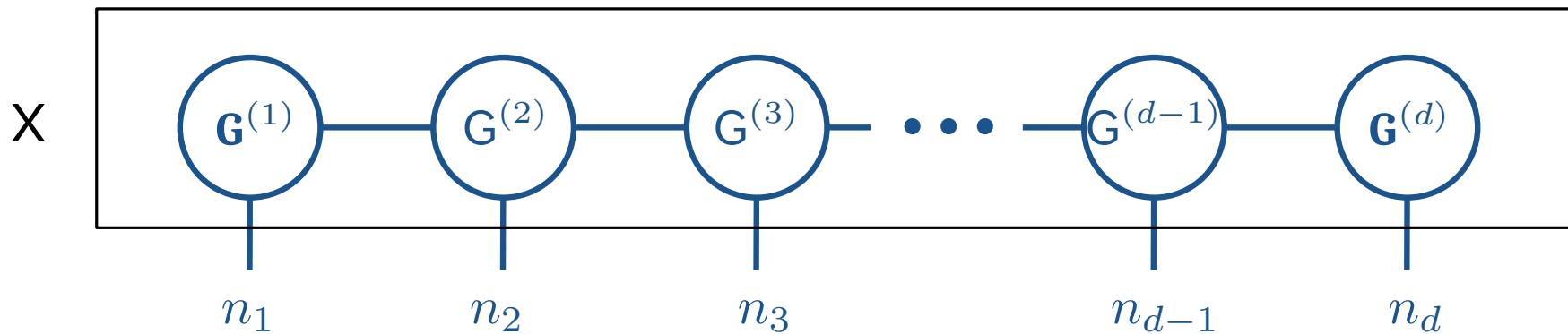
The TT is **non-unique**: for any nonsingular $\mathbf{S} \in \mathbb{R}^{R_i \times R_i}$,

$$\mathbf{G}^{(i)}(n_i) \mathbf{G}^{(i+1)}(n_{i+1}) = \mathbf{G}^{(i)}(n_i) \mathbf{S} \mathbf{S}^{-1} \mathbf{G}^{(i+1)}(n_{i+1}) = \tilde{\mathbf{G}}^{(i)}(n_i) \tilde{\mathbf{G}}^{(i+1)}(n_{i+1})$$

1: Hackbusch, 2012, 2: Vidal, 2003

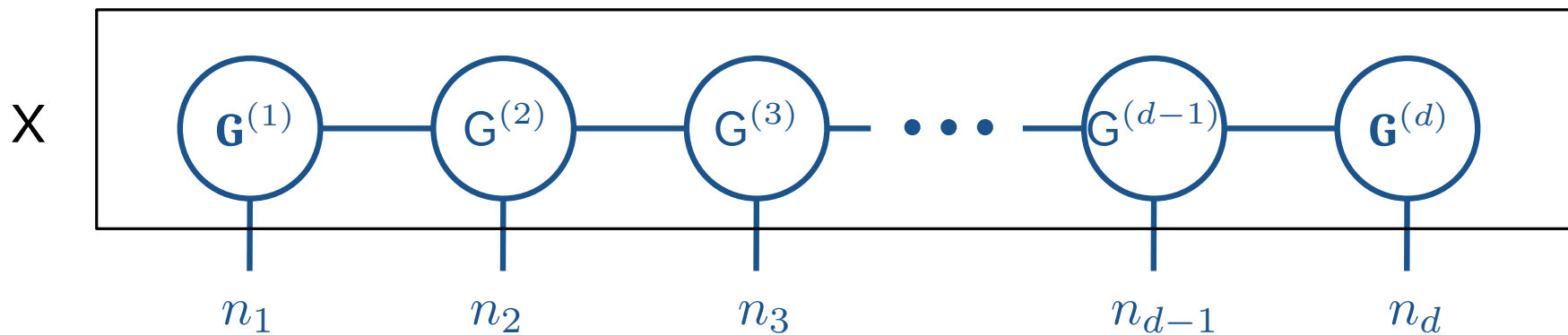
Computation of approximate TT

Standard (algebraic) algorithm: TT-SVD,¹ which performs a sequence of reshapings and truncated SVD, given target TT-ranks (R_1, \dots, R_{d-1}) .



Computation of approximate TT

Standard (algebraic) algorithm: TT-SVD,¹ which performs a sequence of reshapings and truncated SVD, given target TT-ranks (R_1, \dots, R_{d-1}) .



It is quasi-optimal, similarly to the THOSVD:

Thm (Oseledets, 2011): The TT-SVD algorithm satisfies

$$\|X - X_{\text{TT-SVD}}\|_F^2 \leq (d - 1) \|X - X_{\text{best-TT-}(R_1, \dots, R_{d-1})}\|_F^2$$

1: Oseledets, 2011