

Tensor Decompositions and their Applications in Machine Learning

Seminar: Efficient Inference and Large-Scale Machine Learning

Stephan Rabanser

Technische Universität München
Department of Informatics
Data Mining and Analytics
<http://kdd.in.tum.de>

May 7, 2017

Table of contents

Matrix Decomposition: A Motivating Example

Tensor Basics

Tensor Order, Tensor Slices / Fibers, Tensor Rank

Tensor Applications in Machine Learning

Case Study: Learning Mixture of Gaussians

Tensor Decomposition Algorithms

Canonical Polyadic Decomposition (CPD)

Tucker Decomposition

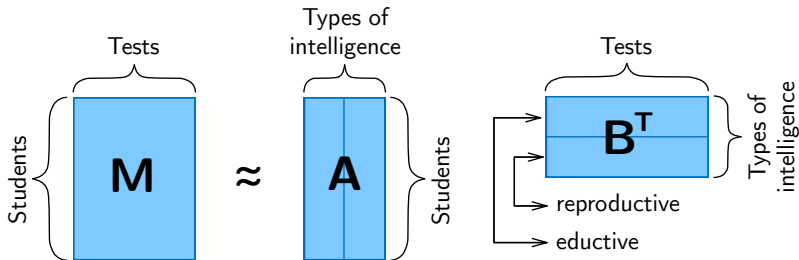
Backup

Matrix Decomposition: A Motivating Example

Spearman's Hypothesis

Charles Spearman, a British psychologist, supposed that human intelligence can be broken down into two factors: **eductive**¹ and **reproductive**² intelligence.

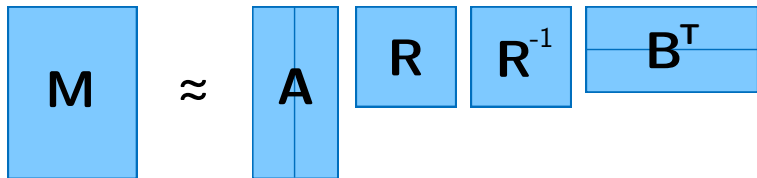
He invented what we refer to as **factor analysis**!



Can we recover **A** and **B** uniquely from **M**? In general, no!

¹ eductive = the ability to make sense out of complexity

² reproductive = the ability to store and reproduce information



$$\begin{aligned} M &= AB^T = ARR^{-1}B^T \\ &= (AR)(R^{-1}B^T) = (AR)(BR^{-T})^T = \tilde{A}\tilde{B}^T \end{aligned}$$

A and B are not determined uniquely unless we impose additional conditions on them (e.g. SVD relies on orthogonality constraints).

Can we achieve weaker uniqueness conditions when moving to higher-dimensional structures (= tensors)? Yes!

Tensor Basics

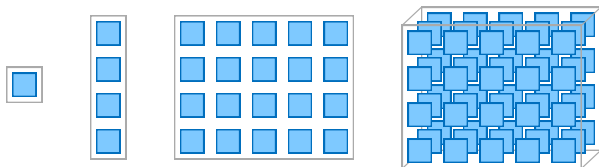
Tensors and Tensor Order

Tensors are generalizations of matrices to higher dimensions and can therefore be seen as multidimensional fields.

The *order/way/mode* of a tensor is the number of its dimensions.

- Zeroth order tensors (scalars): $x \in \mathbb{R}$
- First order tensors (vectors): $\mathbf{x} \in \mathbb{R}^l$
- Second order tensors (matrices): $\mathbf{X} \in \mathbb{R}^{l_1 \times l_2}$
- Third or higher order tensors: $\mathcal{X} \in \mathbb{R}^{l_1 \times l_2 \times \dots \times l_N}$

We can depict zeroth order tensors to third order tensors as follows:

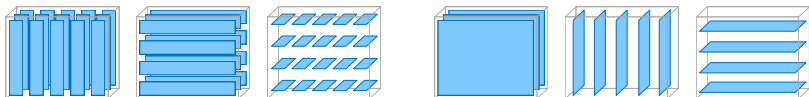


Tensor Fibers and Slices

Fibers are created when fixing all but one index, *slices* (or *slabs*) are created when fixing all but two indices.

For a third order tensor

- fibers are given as $\mathbf{x}_{:jk}$ ($= \mathbf{x}_{jk}$) (column), $\mathbf{x}_{i:k}$ (row), and $\mathbf{x}_{ij:}$ (tube)
- slices are given as $\mathbf{X}_{::k}$ ($= \mathbf{X}_k$) (frontal), $\mathbf{X}_{:j}$ (lateral), and $\mathbf{X}_{i::}$ (horizontal)



Outer Product (Tensor Product)

The outer vector of two n -sized vectors \mathbf{a} , \mathbf{b} is defined as follows:

$$\mathbf{a} \odot \mathbf{b} = \mathbf{a}\mathbf{b}^T$$

The outer product between two vectors produces a matrix.

Inner Product

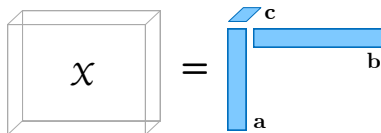
The inner product of two n -sized vectors \mathbf{a} , \mathbf{b} is defined as follows:

$$\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n$$

The inner product between two vectors produces a scalar.

Rank-One Tensors

A N -way tensor is of *rank-one* if it can be strictly decomposed into the (vector) outer product of N vectors.



$$\mathcal{X} = \mathbf{a}^{(1)} \otimes \mathbf{a}^{(2)} \otimes \dots \otimes \mathbf{a}^{(N)} \text{ with } x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} a_{i_2}^{(2)} \dots a_{i_N}^{(N)}$$

Examples

- Rank-one matrix: $\mathbf{X} = \mathbf{a} \otimes \mathbf{b}$
- Rank-one 3-way tensor: $\mathcal{X} = \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$

Tensor Rank

The *rank* of a tensor $\text{rank}(\mathcal{X}) = R$ is defined as the minimum number of rank-one tensors which are needed to produce \mathcal{X} .

$$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r^{(1)} \odot \mathbf{a}_r^{(2)} \odot \dots \odot \mathbf{a}_r^{(n)} = \llbracket \boldsymbol{\lambda}; \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(n)} \rrbracket$$

A specific $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_R]$ is called *factor matrix* and holds the combination of the vectors from the rank-one components as columns. The additional λ_r term is used to absorb the respective weights during normalization of the factor matrices' columns.

Examples

- Rank- R matrix: $\mathbf{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \odot \mathbf{b}_r = \llbracket \boldsymbol{\lambda}; \mathbf{A}, \mathbf{B} \rrbracket$
- Rank- R 3-way tensor: $\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \odot \mathbf{b}_r \odot \mathbf{c}_r$

Tensor Rigidity

Slices of a low-rank tensor are themselves low-rank matrices. A low-rank tensor is not just a collection of low-rank matrices, but there exist relationships between these slices. In fact, the slices are different scalings of the same set of rank-1 matrices.

$$\mathbf{x}_k = \sum_{r=1}^R \left(\mathbf{a}_r \odot \mathbf{b}_r \right) c_{kr}$$

For true underlying factors, subtracting off scalings of the same rank-1 matrix has to decrease the rank of each slice.

$$\mathbf{x}_k - c_k(\mathbf{a} \odot \mathbf{b})$$

This makes tensors way more rigid than matrices!

Multilinear Tensor Transformation

A tensor \mathcal{X} can be transformed on multiple dimensions by hitting each of the dimensions with a (potentially different) transformation tensor from the left side.

$$\tilde{\mathcal{X}} = \mathcal{X}(\mathcal{B}, \mathcal{C}, \mathcal{D}) = \sum_{r=1}^R \lambda_r \langle \mathcal{B} a_r \rangle \langle \mathcal{C} b_r \rangle \langle \mathcal{D} c_r \rangle$$

Hitting a tensor \mathcal{X} with another tensor $\hat{\mathcal{X}}$ on a specific dimension corresponds to an inner product between $\hat{\mathcal{X}}$ and the vector on the corresponding dimension.

Tensor Applications in Machine Learning

- Unsupervised learning settings
 - Topic and mixture modelling
 - Dimensionality reduction
 - Source separation and clustering
- More natural representation for some data
 - Timed data as tensor slices
- Pattern and anomaly detection in graphs
 - Intrusion detection
 - Social network analysis

Typical problems in unsupervised learning settings

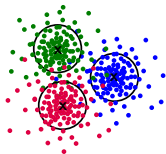
- **Conditions for Identifiability:** We have to determine whether a proposed model contains all relevant information for parameter derivation.
- **Efficient Learning of Latent Variable Models:** We have to determine whether we can estimate the parameters for the proposed model in an efficient (and hopefully optimal) way.

The following approach

- allows us to extract all important information for model construction
- enables us to learn gaussian mixture models efficiently

Gaussian Mixture Model notation

- k : number of mixture components (latent factors).
- $\mathbf{h} \in [\mathbf{e}_1, \dots, \mathbf{e}_k]$: categorical hidden variable represented in basis vector form (= hot-one encoding). $\mathbb{E}[\mathbf{h}] = \mathbf{w}$.
- $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_k]$: means of the mixture components.
- $\boldsymbol{\mu}_C = \mathbf{0}$ (zero mean) and $\boldsymbol{\Sigma}_C = \sigma^2 \mathbf{I}$ (spherical covariance).
- Sample generation: $\mathbf{x} = \mathbf{A}\mathbf{h} + \mathbf{z}$ with $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_C, \boldsymbol{\Sigma}_C)$.



Method of Moments

- Alternative to maximum likelihood estimation.
- Extract model information by building up higher order moments of the underlying probability distribution.

In the GMM case, we can compute the first three moments as follows:

- $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] = \mathbf{A}\mathbf{w} = \sum_{i \in [k]} w_i \mathbf{a}_i$
- $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{x} \odot \mathbf{x}] = \dots = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i + \sigma^2 \mathbf{I}$
- $\boldsymbol{\mathcal{T}} = \mathbb{E}[\mathbf{x} \odot \mathbf{x} \odot \mathbf{x}] = \dots = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i \odot \mathbf{a}_i + \sigma^2 \sum_{i \in [k]} (\boldsymbol{\mu} \odot \mathbf{e}_i \odot \mathbf{e}_i + \mathbf{e}_i \odot \boldsymbol{\mu} \odot \mathbf{e}_i + \mathbf{e}_i \odot \mathbf{e}_i \odot \boldsymbol{\mu})$

Learning Mixture of Gaussians (cont.)

We observe a similar structure in these moments:

- $\boldsymbol{\mu} = \sum_{i \in [k]} w_i \mathbf{a}_i$
- $\boldsymbol{\Sigma} = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i + \sigma^2 \mathbf{I}$
- $\boldsymbol{\mathcal{T}} = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i \odot \mathbf{a}_i + \sigma^2 \sum_{i \in [d]} (\boldsymbol{\mu} \odot \mathbf{e}_i \odot \mathbf{e}_i + \mathbf{e}_i \odot \boldsymbol{\mu} \odot \mathbf{e}_i + \mathbf{e}_i \odot \mathbf{e}_i \odot \boldsymbol{\mu})$

σ^2 is estimated by the smallest eigenvalue of the second moment $\boldsymbol{\Sigma}$. We can therefore reduce the moment computation to:

- $\mathbf{m}_1 = \sum_{i \in [k]} w_i \mathbf{a}_i$
- $\mathbf{M}_2 = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i$
- $\mathbf{M}_3 = \sum_{i \in [k]} w_i \mathbf{a}_i \odot \mathbf{a}_i \odot \mathbf{a}_i$

Many latent variable models can be reduced to this form!

If we now assume \mathbf{A} to be (columnwise) orthogonal, then the \mathbf{a}_i 's are the eigenvectors of the third-order moment tensor \mathcal{M}_3 . If we knew for example \mathbf{a}_1 , then hitting \mathcal{M}_3 with \mathbf{a}_1 on two dimensions yields:

$$\mathcal{M}_3(\mathbf{I}, \mathbf{a}_1, \mathbf{a}_1) = \sum_{i \in [k]} w_i \langle \mathbf{a}_i, \mathbf{a}_1 \rangle^2 \mathbf{a}_i = w_1 \mathbf{a}_1$$

This directly corresponds to the concept of matrix eigenvectors! However, we are still facing two problems at this point:

- Can we transform any given tensor to an orthogonal tensor?
- How do we find the eigenvectors of an orthogonal tensor?

Orthogonalization Through Whitening

Using the second order moment, we can obtain a whitening transformation, that orthogonalizes \mathbf{A} . The transformation has to be invertible, which is why we want \mathbf{A} to have full column rank.

$$\mathbf{M}_2 = \mathbf{U}\text{Diag}(\tilde{\lambda})\mathbf{U}^T \Rightarrow \mathbf{W} = \mathbf{U}\text{Diag}(\tilde{\lambda}^{-1/2})$$

The orthogonalized moment \mathcal{V} can then be computed as follows:

$$\begin{aligned}\mathcal{V} = \mathfrak{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) &= \sum_{i \in [k]} \frac{1}{\sqrt{w_i}} (\mathbf{W}^T \mathbf{a}_i \sqrt{w_i})^{\odot 3} \\ &= \sum_{i \in [k]} \lambda_i \mathbf{v}_i \odot \mathbf{v}_i \odot \mathbf{v}_i\end{aligned}$$

Learning Mixture of Gaussians (cont.)

Tensor Decomposition Through Power Method

Now we can find the eigenvectors of the orthogonalized tensor \mathcal{V} by applying the tensor power method, a natural extension of the matrix power method. After randomly initializing a vector \mathbf{v} we can feed the vector into a tensor power step and repeat this step until convergence:

$$\mathbf{v} \mapsto \frac{\mathcal{V}(\mathbf{I}, \mathbf{v}, \mathbf{v})}{\|\mathcal{V}(\mathbf{I}, \mathbf{v}, \mathbf{v})\|}$$

Afterwards we deflate the tensor, i.e. remove eigenvalue λ with eigenvector \mathbf{v} that we just extracted from the tensor as follows:

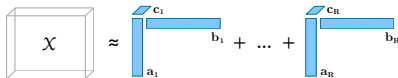
$$\mathcal{V}_{\text{next}} = \mathcal{V} - \lambda \mathbf{v} \odot \mathbf{v} \odot \mathbf{v}$$

This process can be repeated until we have extracted the k dominant eigenvalues/-vectors.

Tensor Decomposition Algorithms

Canonical Polyadic Decomposition (CPD)

The Canonical Polyadic Decomposition (CPD) aims at approximating a tensor \mathcal{X} as good as possible with a low-rank tensor $\hat{\mathcal{X}}$.



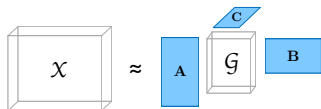
$$\min_{\hat{\mathcal{X}}} \|\mathcal{X} - \hat{\mathcal{X}}\| \quad \text{with} \quad \hat{\mathcal{X}} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \odot \mathbf{b}_r \odot \mathbf{c}_r = [\![\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}]\!]$$

Compute \mathbf{A} , \mathbf{B} , and \mathbf{C} through:

- **Jennrich's algorithm:** eigendecomposition of tensor slices
- **Alternating least squares (ALS) algorithm:** fix all factor matrices except for one in order to optimize for the non-fixed matrix

Tucker Decomposition

The Tucker decomposition can be thought of as a higher-order PCA. $\mathcal{G} \in \mathbb{R}^{P \times Q \times R}$ is called core tensor.



$$\min_{\hat{\mathcal{X}}} \|\mathcal{X} - \hat{\mathcal{X}}\| \quad \text{with} \quad \hat{\mathcal{X}} = \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R g_{pqr} \mathbf{a}_r \odot \mathbf{b}_r \odot \mathbf{c}_r$$
$$= \llbracket \mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$$

Compute \mathbf{A} , \mathbf{B} , and \mathbf{C} through:

- Higher order singular value decomposition (HOSVD)
- Higher order orthogonal iteration (HOOI)

Thanks! :)

Backup

Tensor Reorderings

We can *vectorize* a given matrix $\mathbf{X} \in \mathbb{R}^{I \times J}$ by vertically stacking the columns of \mathbf{X} into a tall vector.

$$\text{vec}(\mathbf{X}) = \begin{bmatrix} \mathbf{x}_{:1} \\ \mathbf{x}_{:2} \\ \vdots \\ \mathbf{x}_{:J} \end{bmatrix}$$

The mode- n matricization of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ (denoted $\mathbf{X}_{(n)}$) turns the mode- n fibers of \mathcal{X} into the columns of $\mathbf{X}_{(n)}$. Let $x \in \mathcal{X}$ be an element of a tensor and $m \in \mathbf{M}$ be an element of the unfolded tensor. Then the mode- n matricization is defined as:

$$x_{i_1, i_2, \dots, i_N} \mapsto m_{i_n, j} \text{ with } j = 1 + \sum_{\substack{k=1 \\ k \neq n}}^N \left((i_k - 1) \prod_{\substack{m=1 \\ m \neq n}}^{k-1} I_m \right)$$

Example

Let \mathcal{X} be a tensor with the following frontal slices:

$$\mathbf{X}_1 = \begin{bmatrix} 1 & 4 & 7 & 10 \\ 2 & 5 & 8 & 11 \\ 3 & 6 & 9 & 12 \end{bmatrix} \quad \mathbf{X}_2 = \begin{bmatrix} 13 & 16 & 19 & 22 \\ 14 & 17 & 20 & 23 \\ 15 & 18 & 21 & 24 \end{bmatrix}$$

Then the three mode- n matricizations are:

$$\mathbf{X}_{(1)} = \begin{bmatrix} 1 & 4 & 7 & 10 & 13 & 16 & 19 & 22 \\ 2 & 5 & 8 & 11 & 14 & 17 & 20 & 23 \\ 3 & 6 & 9 & 12 & 15 & 18 & 21 & 24 \end{bmatrix} \quad \mathbf{X}_{(2)} = \begin{bmatrix} 1 & 2 & 3 & 13 & 14 & 15 \\ 4 & 5 & 6 & 16 & 17 & 18 \\ 7 & 8 & 9 & 19 & 20 & 21 \\ 10 & 11 & 12 & 22 & 23 & 24 \end{bmatrix}$$

$$\mathbf{X}_{(3)} = \begin{bmatrix} 1 & 2 & 3 & 4 & \cdots & 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 & \cdots & 21 & 22 & 23 & 24 \end{bmatrix}$$

Kronecker Product

The Kronecker product between two arbitrarily-sized matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{K \times L}$ is a generalization of the outer product from vectors to matrices.

$$\begin{aligned} \mathbf{A} \otimes \mathbf{B} &:= \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1J}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2J}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}\mathbf{B} & a_{I2}\mathbf{B} & \cdots & a_{IJ}\mathbf{B} \end{bmatrix} \\ &= [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \mathbf{a}_1 \otimes \mathbf{b}_2 \quad \cdots \quad \mathbf{a}_J \otimes \mathbf{b}_{L-1} \quad \mathbf{a}_J \otimes \mathbf{b}_L] \end{aligned}$$

The resulting matrix will be of size $(IK) \times (JL)$.

Khatri-Rao Product

The Khatri-Rao product between two matrices $\mathbf{A} \in \mathbb{R}^{I \times K}$ and $\mathbf{B} \in \mathbb{R}^{J \times K}$ corresponds to the column-wise Kronecker product.

$$\mathbf{A} \odot \mathbf{B} := [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \mathbf{a}_2 \otimes \mathbf{b}_2 \quad \cdots \quad \mathbf{a}_K \otimes \mathbf{b}_K]$$

The resulting matrix will be of size $(IJ) \times K$.

Hadamard Product

The Hadamard product between two same-sized matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{I \times J}$ corresponds to the element-wise matrix product.

$$\mathbf{A} * \mathbf{B} := \begin{bmatrix} a_{11}b_{11} & a_{12}b_{12} & \cdots & a_{1J}b_{1J} \\ a_{21}b_{21} & a_{22}b_{22} & \cdots & a_{2J}b_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}b_{I1} & a_{I2}b_{I2} & \cdots & a_{IJ}b_{IJ} \end{bmatrix}$$

The resulting matrix will also be of size $I \times J$.

Important Matrix/Tensor Products (cont.)

n -mode Product

The n -mode product of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ and a matrix $\mathbf{M} \in \mathbb{R}^{J \times I_n}$ is denoted by $\mathcal{Y} = \mathcal{X} \times_n \mathbf{M}$. Intuitively, each mode- n fiber is multiplied by the matrix \mathbf{M} .

$$(\mathcal{X} \times_n \mathbf{M})_{i_1 \dots i_{n-1} j i_{n+1} \dots i_N} = \sum_{i_n=1}^{I_n} x_{i_1 \dots i_N} m_{j i_n}$$

The n -mode product can also be expressed through matricized tensors as $\mathbf{Y}_{(n)} = \mathbf{M} \mathbf{X}_{(n)}$.

For vectors, we analogously define:

$$(\mathcal{X} \times_n \mathbf{v})_{i_1 \dots i_{n-1} i_{n+1} \dots i_N} = \sum_{i_n=1}^{I_n} x_{i_1 \dots i_N} v_{i_n}$$

Jennrich's Algorithm

If \mathbf{A} , \mathbf{B} , and \mathbf{C} are full rank, then $\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \odot \mathbf{b}_r \odot \mathbf{c}_r$ is unique up to rank permutation and feature scaling.

We can use Jennrich's algorithm to recover the factor matrices.

1. Choose random vectors \mathbf{x} and \mathbf{y} .
2. $\mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{x}) = \sum_{r=1}^R \langle \mathbf{c}_r, \mathbf{x} \rangle \mathbf{a}_r \odot \mathbf{b}_r = \mathbf{A} \text{Diag}(\langle \mathbf{c}_r, \mathbf{x} \rangle) \mathbf{B}^T$.
3. $\mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{y}) = \sum_{r=1}^R \langle \mathbf{c}_r, \mathbf{y} \rangle \mathbf{a}_r \odot \mathbf{b}_r = \mathbf{A} \text{Diag}(\langle \mathbf{c}_r, \mathbf{y} \rangle) \mathbf{B}^T$.
4. $\mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{x}) \mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{y})^\dagger = \mathbf{A} \text{Diag}(\langle \mathbf{c}_r, \mathbf{x} \rangle) \text{Diag}(\langle \mathbf{c}_r, \mathbf{y} \rangle)^\dagger \mathbf{A}^\dagger$
 \Rightarrow eigendecomposition to find \mathbf{A}
5. $\mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{x})^\dagger \mathcal{X}(\mathbf{I}, \mathbf{I}, \mathbf{y}) = (\mathbf{B}^T)^\dagger \text{Diag}(\langle \mathbf{c}_r, \mathbf{x} \rangle)^\dagger \text{Diag}(\langle \mathbf{c}_r, \mathbf{y} \rangle) \mathbf{B}^T$
 \Rightarrow eigendecomposition to find \mathbf{B}
6. Match up the factors and solve a linear system to find \mathbf{C} .

However, we are still not using the full tensor structure.

Alternating Least Squares (ALS)

Fix all factor matrices except for one in order to optimize for the non-fixed matrix and then repeat this step for every matrix repeatedly until some stopping criterion is satisfied.

$$\mathbf{A} \leftarrow \arg \min_{\mathbf{A}} \|\mathbf{X}_{(1)} - (\mathbf{C} \odot \mathbf{B})\mathbf{A}^T\|$$

$$\mathbf{B} \leftarrow \arg \min_{\mathbf{B}} \|\mathbf{X}_{(2)} - (\mathbf{C} \odot \mathbf{A})\mathbf{B}^T\|$$

$$\mathbf{C} \leftarrow \arg \min_{\mathbf{C}} \|\mathbf{X}_{(3)} - (\mathbf{B} \odot \mathbf{A})\mathbf{C}^T\|$$

The factor matrices are unique if the following holds:

- $\min_{1, \dots, N} \left(\prod_{\substack{m=1 \\ m \neq n}}^N \text{rank}(\mathbf{A}^{(m)}) \right) \geq R$
- For $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$: $R \leq K$ and $R(R-1) \leq \frac{I(I-1)J(J-1)}{2}$

Optimal ALS Solution

$$\hat{\mathbf{A}} = \mathbf{X}_{(1)}[(\mathbf{C} \odot \mathbf{B})^T]^\dagger = \mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B})^\dagger$$

$$\hat{\mathbf{B}} = \mathbf{X}_{(2)}[(\mathbf{C} \odot \mathbf{A})^T]^\dagger = \mathbf{X}_{(2)}(\mathbf{C} \odot \mathbf{A})(\mathbf{C}^T \mathbf{C} * \mathbf{A}^T \mathbf{A})^\dagger$$

$$\hat{\mathbf{C}} = \mathbf{X}_{(3)}[(\mathbf{B} \odot \mathbf{A})^T]^\dagger = \mathbf{X}_{(3)}(\mathbf{B} \odot \mathbf{A})(\mathbf{B}^T \mathbf{B} * \mathbf{A}^T \mathbf{A})^\dagger$$

However, this setting can be numerically unstable due to potential ill-conditioning.

CPD Through ALS

```
procedure CP-ALS( $\mathbf{X}$ ,  $R$ )  
  initialize  $\mathbf{A}^{(n)} \in R^{I_n \times R}$  for  $n = 1, \dots, N$   
  repeat  
    for  $n = 1, \dots, N$  do  
       $\mathbf{V} \leftarrow \mathbf{A}^{(1)T} \mathbf{A}^{(1)} * \dots * \mathbf{A}^{(n-1)T} \mathbf{A}^{(n-1)} * \mathbf{A}^{(n+1)T} \mathbf{A}^{(n+1)} *$   
       $\hookrightarrow \dots * \mathbf{A}^{(N)T} \mathbf{A}^{(N)}$   
       $\mathbf{A}^{(n)} \leftarrow \mathbf{X}_{(n)}(\mathbf{A}^{(N)} \odot \dots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \dots \odot \mathbf{A}^{(1)}) \mathbf{V}^\dagger$   
      normalize columns of  $\mathbf{A}^{(n)}$  (optional)  
      store norms as  $\lambda$  (optional)  
    end for  
  until stopping criterion satisfied  
  return  $\lambda, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}$   
end procedure
```

Tensor Rank Peculiarities

- There is no trivial algorithm to determine the rank of a tensor as the problem is NP-hard.
- The rank of a tensor with real entries can be different over \mathbb{R} and \mathbb{C} .
- The maximum rank of a tensor represents the largest attainable rank. The typical rank is any rank that occurs with probability ≥ 0 . These two ranks are usually not the same for tensors (while they are for matrices).
- The rank of a tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ is constant under permutation and bounded from above by the following equation:

$$\text{rank}(\mathcal{X}) \leq \min\{IJ, IK, JK\}$$

Tucker Decomposition (cont.)

The Tucker decomposition is not unique! This follows from the fact that all columns of the factor matrices can interact with arbitrary other columns through the core tensor.

CPD can be thought of a special case of the Tucker decomposition with a superdiagonal core.

Tucker Decomposition (cont.)

n-rank

The n-rank of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$, denoted $R_n = \text{rank}_n(\mathcal{X})$, corresponds to the column rank of $\mathbf{X}_{(n)}$.

Higher Order Singular Value Decomposition (HOSVD)

```
procedure HOSVD( $\mathcal{X}, R_1, \dots, R_N$ )  
  for  $n = 1, \dots, N$  do  
     $\mathbf{A}^{(n)} \leftarrow R_n$  leading left singular vectors of  $\mathbf{X}_{(n)}$   
  end for  
   $\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^{(1)T} \times_2 \dots \times_N \mathbf{A}^{(N)T}$   
  return  $\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}$   
end procedure
```

This algorithm does not give an optimal fit, but it is a good starting point for a successive application of the ALS algorithm!

Higher Order Orthogonal Iteration (HOOI)

```
procedure HOOI( $\mathcal{X}, R_1, \dots, R_N$ )  
  initialize  $\mathbf{A}^{(n)} \in R^{I_n \times R}$  for  $n = 1, \dots, N$  using HOSVD  
  repeat  
    for  $n = 1, \dots, N$  do  
       $\mathcal{Y} \leftarrow \mathcal{X} \times_1 \mathbf{A}^{(1)T} \times_2 \dots \times_{n-1} \mathbf{A}^{(n-1)T} \times_{n+1} \mathbf{A}^{(n+1)T} \times_{n+2}$   
       $\hookrightarrow \dots \times_N \mathbf{A}^{(N)T}$   
       $\mathbf{A}^{(n)} \leftarrow R_n$  leading left singular vectors of  $\mathcal{Y}_{(n)}$   
    end for  
  until stopping criterion satisfied  
   $\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{A}^{(1)T} \times_2 \dots \times_N \mathbf{A}^{(N)T}$   
  return  $\mathcal{G}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}$   
end procedure
```

CPD vs. Tucker

CPD	Tucker
linear tensor decomposition	multilinear tensor decomposition
column-matching interaction	various interactions
true underlying factors	compression
unique	not unique

- **Individual differences in scaling (INDSCAL)**: special case of CPD for 3-way tensors with two symmetric modes.
- **Parallel factors for cross products (PARAFAC2)**: CPD which can be applied to a collection of matrices with the same amount of columns, but a different number of rows.
- **CANDECOMP with linear constraints (CANDELINC)**: add constraints to CPD's factor matrices through domain/user knowledge.
- **Decomposition into directional components (DEDICOM)**
- **PARAFAC and Tucker 2 (PARATUCK2)**