

Mathematics, Pusan National University

LINEAR ALGEBRA AND LEARNING FROM DATA

Factoring Matrices and Tensors : Positive and Sparse

Taehyeong Kim
th_kim@pusan.ac.kr

August 28, 2020



Introduction

Nonnegative Matrix Factorization (NMF)

Text Mining

Facial Feature Extraction

Optimality Conditions for Nonnegative U and V

Sparse Principal Components

Tensors

Example 1 : A Color Image is a Tensor with 3 Slices

Example 2 : The Derivative $\partial \mathbf{w} / \partial A$ of $\mathbf{w} = A\mathbf{v}$

Example 3 : The Joint Probability Tensor

Tensor Multiplications

The Norm and Rank of a Tensor

The CP Decomposition of a Tensor

Matricized Form of a Tensor \mathcal{I}

The Khatri-Rao Product $A \odot B$

Computing the CP Decomposition of \mathcal{I}

The Tucker Decomposition

$$A \geq 0 \quad U \geq 0 \quad V \geq 0$$

$$A = UV$$



Here are factorizations of A and T with new and important properties :

Properties

★ **Nonnegative Matrices**

★ **Sparse and Nonnegative**

★ **CP Tensor Decomposition**

$$\begin{aligned} & \min \|A - UV\|_F^2 \\ & \text{with } U \geq 0 \text{ and } V \geq 0 \\ & \min \|A - UV\|_F^2 + \lambda \|UV\|_N \\ & \text{with } U \geq 0 \text{ and } V \geq 0 \\ & \min \|T - \sum_{i=1}^R \mathbf{a}_i \otimes \mathbf{b}_i \otimes \mathbf{c}_i\| \end{aligned}$$

To compute a factorization $A = UV$, we introduce a simple **alternating iteration.**
Update U with V fixed, then update V with U fixed.



The goal of NMF

Approximating a nonnegative matrix $A \geq 0$ by a lower rank product UV of two nonnegative matrices $U \geq 0$ and $V \geq 0$.

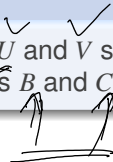
NMF and SPCA

NMF

Find ~~nonnegative~~ matrices U and V so that $A \approx \underline{UV}$

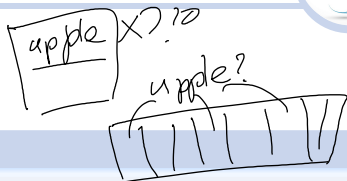
SPCA

Find sparse low rank matrices B and C so that $A \approx \underline{BC}$





We can use NMF to discover new characteristics of data.



TF-IDF

TF(Term Frequency) : The weight of a term that occurs in a document is simply proportional to the term frequency.

DF(Document Frequency) : The specificity of a term can be quantified as an inverse function of the number of documents in which it occurs.

TF-IDF (Term Frequency Inverse Document Frequency) :

$$\frac{TF}{DF}$$



Example

Suppose that the matrix V made up of book titles and words created using TF-IDF is as follows :

	협상	스타트업	투자	비즈니스	데이터
$V =$ 협상의 법칙	0.9	0	0.3	0.8	0
스타트업	0	0.8	0.7	0.9	0.3
빅데이터	0	0	0.5	0	0.8



Example

Then the matrices W and H such that $\underline{WH} \approx V$ are

es W and H such that $|WH| \approx V$ are

		특징1	특징2	특징3	
$\underline{W} =$	협상의 법칙	0	1.2410	0	≥ 0
	스타트업	1.4248	0	0	
	빅데이터	0	0	0.9434	

		협상	스타트업	투자	비즈니스	데이터	
$\underline{H} =$	특징1	0	0.5615	0.4913	0.6317	0.2106	≥ 0
	특징2	0.7252	0	0.2417	0.6447	0	
	특징3	0	0	0.5300	0	0.8480	

W is called the weight matrix and H is called the feature matrix.



For example, there are 20 $\overbrace{45 \times 40 \text{ face}}^{\times 1}$ images like this:



Figure: 45×40 face images

If we make a vector by connecting the pixel brightness values in an image, we can think of each face image as a $45 \times 40 = 1800$ dimensional vector. Now, if we perform PCA with these 1800 dimensional point data, we can obtain the same number of principal component vectors as the number of dimensions of the data. **Eigenface** is the reinterpretation of the principal component vectors thus obtained as images.

Facial Feature Extraction



Figure: reconstruction using k eigenfaces

If we use a large number of eigenfaces, we can see approximate results that are almost similar to those of the original face, but as k decreases, the unique facial features of the individual disappear and the common facial features remain.

It can be interpreted in various ways as dimension reduction, data compression, and noise reduction.



By using this, we can detect face.

1. After collecting many face samples, get their eigenfaces and select only the first k eigenfaces.
2. When the test image x comes in, calculate how close it is to the original image x when x is reconstruction using only k eigenfaces.
3. If x is perfectly approximated by combining k eigenfaces, x is very likely to be a face. ←
4. Another criterion of judgement is how close x_k is to the average face. ←

Optimality Conditions for Nonnegative U and V



$$(UV - A)V^T$$

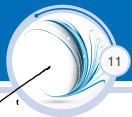
Given $A \geq 0$, here are the conditions for $U \geq 0$ and $V \geq 0$ to minimize $\|A - UV\|_F^2$

$$\begin{cases} Y = UVV^T - AV^T \geq 0 & \text{with } Y_{ij} \text{ or } U_{ij} = 0 \text{ for all } i, j \\ Z = U^T UV - U^T A \geq 0 & \text{with } Z_{ij} \text{ or } V_{ij} = 0 \text{ for all } i, j \end{cases}$$

Those last conditions already suggest that U and V may turn out to be sparse. ✓
The calculation will be shown in chapter 3. //

$$\boxed{UV - A} \quad \boxed{V}^T$$

Sparse Principal Components

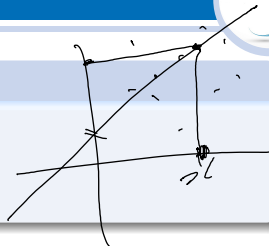


Multiple Linear Regression

$$\arg \min_{\mathbf{w}, \mathbf{b}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \arg \min_{\mathbf{w}, \mathbf{b}} MSE$$

where $\hat{\mathbf{y}} = \mathbf{w}_0 x_1 + \dots + \mathbf{w}_p x_p + \mathbf{b}$

$\mathbb{R} \rightarrow \mathbb{R}$



LASSO

Lasso is simply a linear regression method with an l^1 -norm penalty.

$$\arg \min_{\mathbf{w}, \mathbf{b}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^m |w_j| = \arg \min_{\mathbf{w}, \mathbf{b}} MSE + \text{penalty}$$

α is a parameter that controls the effect of the penalty.
It can create a model that avoids overfitting.



LASSO

$$\text{Minimize } \underbrace{\|A\mathbf{x} - \mathbf{b}\|^2}_{//} + \lambda \underbrace{\sum_{k=1}^n |x_k|}_{//}$$

Elastic net

$$\text{Minimize } \underbrace{\|A\mathbf{x} - \mathbf{b}\|_2^2}_{//} + \lambda \underbrace{\|\mathbf{x}\|_1}_{//} + \beta \underbrace{\|\mathbf{x}\|_2^2}_{//} \quad 3.4$$

Tensors



$$(3, 2, 2, \dots) \quad (a_{i_1 i_2 \dots i_m}) : \text{size} = n_1 \times \dots \times n_m$$

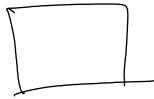
Generally tensor is multi-array such that $\mathcal{A} = (a_{i_1 \dots i_m})$ $a_{i_1 \dots i_m} \in \mathbb{F}$: field where $i_j = 1, \dots, n_j$, for $j = 1, \dots, m$. Let's take an example with a matrix we know well.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = (a_{i_1 i_2})$$

Matrix $A = (a_{i_1 i_2})$ is tensor that $i_1 = 1, 2$, $i_2 = 1, 2, 3$. m is order of tensor, (n_1, \dots, n_m) is dimension of tensor. Similar to the size of matrix A is $2 \times 3 = 6$, the size of tensor \mathcal{A} is $n_1 \times \dots \times n_m$. If $n_1 = \dots = n_m$, \mathcal{A} is m th order n -dimensional tensor.

$$\underline{\underline{T_{m, n}}}$$

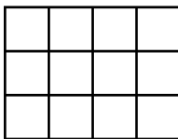
$$\underbrace{i_1 \dots i_m}_{\sim n}$$



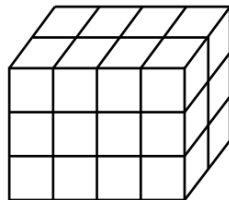
$$\begin{array}{c} a_{i_1 i_2 i_3} \\ \hline 123 \quad 123 \quad 123 \end{array}$$



vector
 x in \mathbb{R}^3



matrix
 A in $\mathbb{R}^{3 \times 4}$

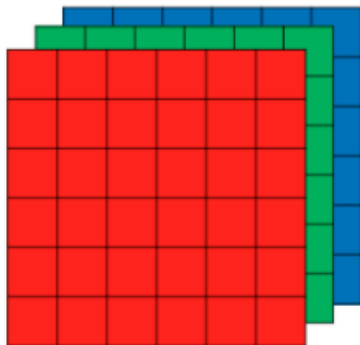


tensor
 T in $\mathbb{R}^{3 \times 4 \times 2}$

Figure: vector, matrix and tensor

order: 3
dim: ~~3x4x2~~
(3, 4, 2)

Example 1 : A Color Image is a Tensor with 3 Slices



6 x 6 x 3

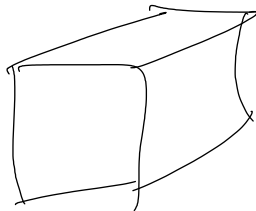


Figure: RGB image data and tensor

Example 2 : The Derivative $\partial \mathbf{w} / \partial A$ of $\mathbf{w} = A\mathbf{v}$



Let A be a **weight** matrix for deep learning. That matrix multiplies a vector \mathbf{v} to produce $\mathbf{w} = A\mathbf{v}$. In matrix multiplication, we know that row j of A has no effect on row i of $\mathbf{w} = A\mathbf{v}$. So the derivative formula includes the symbol δ_{ij} . The derivatives of the linear function $\mathbf{w} = A\mathbf{v}$ with respect to the weights A_{jk} are in \mathcal{T} : \Rightarrow

$$\boxed{\mathcal{T}} = \mathcal{I}$$

$$\mathcal{T}_{ijk} = \frac{\partial w_i}{\partial A_{jk}} = v_k \delta_{ij}$$

$$\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

This tensor \mathcal{T} is interesting:

1. The slices $k = \text{constant}$ are multiples v_k of the identity matrix.
2. The key function of deep learning connects each layer of a neural net to the next layer. If one layer contains a vector \mathbf{v} , the next layer contains the vector $\mathbf{w} = (A\mathbf{v} + \mathbf{b})_+$. A is a matrix of weights. We optimize those weights to match the training data. So the derivatives of the loss function L will be zero for the optimal weights.

Example 3 : The Joint Probability Tensor



$$\begin{array}{l} \frac{10}{N} \quad I = \{8, 9, 10\} \quad H = \{100, 110, 120\} \quad W = \{20, 25, 30\} \end{array}$$

Suppose we measure age a in years and height h and weight w . We put N children into I age groups and J height groups and K weight groups. Then all children will be in the (i, j, k) group. For that random child,

$$\text{Probability of age group } i = \frac{a_i}{N} \quad \text{height group } j = \frac{h_j}{N} \quad \text{weight group } k = \frac{w_k}{N}$$

This naturally leads to **joint probabilities** $p_{ijk} = \frac{\text{number of } i \text{ age, } j \text{ height, } k \text{ weight children}}{N}$.

$$p_{ijk}$$



$$\mathcal{A} = (a_{i_1 \dots i_m})$$

Definition

Tensor Outer Product We use \otimes to denote tensor outer product; that is for any two tensors $\mathcal{A} \in T_{m, \mathbb{Q}}$ and $\mathcal{B} \in T_{p, \mathbb{Q}}$

$$\underline{\mathcal{A} \otimes \mathcal{B}} = (a_{i_1 \dots i_m} b_{i_{m+1} \dots i_{m+p}}) \quad (1)$$

$$\mathcal{A} \otimes \mathcal{B} \in T_{(m+p), n}$$

Tensor Multiplications

$$\mathbb{X} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\mathbb{X} \otimes \mathbb{X} =$$



symmetric rank-one tensor

$$\underline{\underline{\mathbf{x}^{\otimes k}}} \equiv \underbrace{\mathbf{x} \otimes \cdots \otimes \mathbf{x}}_{k \text{ times}} = \underline{\underline{(x_{i_1} \cdots x_{i_k})}} \in T_{k,n} \quad (2)$$

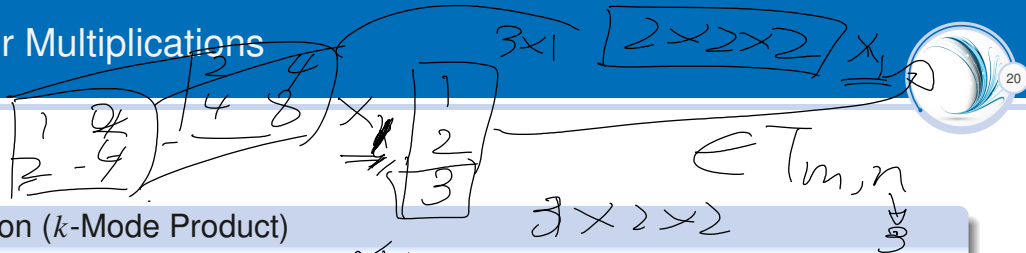
Obviously, $\mathbf{x}^{\otimes k} \in S_{k,n}$, and it called a symmetric rank-one tensor when $\mathbf{x} \neq \mathbf{0}$.

✓ rank-one tensor

More generally, let $\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_n^{(i)})^T \in \mathbb{R}^n$ for $i \in [m]$ and $\alpha \in \mathbb{R}$. Then

$\alpha \mathbf{x}^{(1)} \otimes \mathbf{x}^{(2)} \otimes \cdots \otimes \mathbf{x}^{(m)}$ is a tensor in $T_{m,n}$ with i_1, \dots, i_m th entry as $\alpha x_{i_1}^{(1)} \cdots x_{i_m}^{(m)}$. Such a tensor (not necessarily symmetric) is called a rank-one tensor in $T_{m,n}$.

Tensor Multiplications



Definition (k -Mode Product)

For any $\mathcal{A} \in T_{m,n}$ and any $P = (p_{ij}) \in \mathbb{R}^{p \times n}$, and for any given $k \in [m]$, the k -mode product of \mathcal{A} and P , denoted as $\mathcal{A} \times_k P$, is defined by

$$(\mathcal{A} \times_k P)_{i_1 \dots i_{k-1} j i_{k+1} \dots i_m} = \sum_{i_k=1}^n a_{i_1 \dots i_{k-1} i_k i_{k+1} \dots i_m} p_{i, i_k}, \quad \forall i_l \in [n], l \in [m], l \neq k, \forall j \in [p] \quad (3)$$

By this product, the size of tensor is changed from $n \times \dots \times n$ to $n \times \dots \times p \times \dots \times n$.



linear operator $P^m(\cdot)$

If we do the k-mode product of \mathcal{A} and P for all possible $k \in [n]$ as

$$P^m(\mathcal{A}) = \mathcal{A} \times_1 P \times_2 \cdots \times_m P$$

More specifically,

$$P^m(\mathcal{A}) = \left(\sum_{i_1, \dots, i_m=1}^n a_{i_1, \dots, i_m} p_{j_1 i_1} \cdots p_{j_m i_m} \right) \in T_{m,p},$$

$$\forall \mathcal{A} = (a_{i_1, \dots, i_m}) \in T_{m,n}.$$

(4)

For $\mathbf{x}^T = (x_1, \dots, x_n)$, the following frequently used notations are given as below:

$$\left\| \mathcal{A} \mathbf{x}^{m-2} \right\| \equiv \mathcal{A} \times_3 \mathbf{x}^T \times_4 \cdots \times_m \mathbf{x}^T = \left(\sum_{i_3, \dots, i_m=1}^n a_{i j i_3 \dots i_m} x_{i_3} \cdots x_{i_m} \right) \in \mathbb{R}^{n \times n} \quad (5)$$

$$\left\| \mathcal{A} \mathbf{x}^{m-1} \right\| \equiv \mathcal{A} \times_2 \mathbf{x}^T \times_3 \cdots \times_m \mathbf{x}^T = \left(\sum_{i_2, \dots, i_m=1}^n a_{i i_2 \dots i_m} x_{i_2} \cdots x_{i_m} \right) \in \mathbb{R}^n \quad (6)$$

$$\underline{\underline{\mathcal{A} \mathbf{x}^m}} \equiv \mathcal{A} \times_1 \mathbf{x}^T \times_2 \cdots \times_m \mathbf{x}^T = \left(\sum_{i_1, \dots, i_m=1}^n \underbrace{a_{i_1 \dots i_m} x_{i_1} \cdots x_{i_m}} \right) \in \mathbb{R} \quad (7)$$



Definition (Inner Product)

For any two tensor $\mathcal{A} = (a_{i_1 \dots i_m})$, $\mathcal{B} = (b_{i_1 \dots i_m}) \in T_{m,n}$, the inner product of \mathcal{A} and \mathcal{B} , denoted as $\mathcal{A} \bullet \mathcal{B}$, is defined as

$$\mathcal{A} \bullet \mathcal{B} = \sum_{i_1, \dots, i_m=1}^n a_{i_1 \dots i_m} b_{i_1 \dots i_m}. \quad (8)$$

Frobenious norm of \mathcal{A}

$$\|\mathcal{A}\|_F = \sqrt{\mathcal{A} \bullet \mathcal{A}}$$

Tensor Multiplication

Hadamard Product



$$A \cdot B$$

Definition (Hadamard Product)

For any two tensor $\mathcal{A} = (a_{i_1 \dots i_m})$, $\mathcal{B} = (b_{i_1 \dots i_m}) \in T_{m,n}$, the Hadamard product of \mathcal{A} and \mathcal{B} , denoted as $\mathcal{A} \circ \mathcal{B}$, is defined as

$$\mathcal{A} \circ \mathcal{B} = (a_{i_1 \dots i_m} b_{i_1 \dots i_m}) \in T_{m,n} \quad (9)$$



3 order tensor

Consider 3 dimension tensor $\mathcal{T} = (t_{ijk})$. Frobenius norm of a tensor : **Add all \mathcal{T}_{ijk}^2 to find $\|\mathcal{T}\|$** . The theory of tensors is still part of linear algebra (or perhaps multilinear algebra).

Just like a matrix, a tensor can have two different roles in science and engineering :

1. A tensor can multiply vectors, matrices, or tensors. Then it is a **linear operator**.
2. A tensor can **contain data**. Its entries could give the brightness of pixels in an image. A color image is 3-way, stacking RGB. A color video will be a 4-way tensor.

The rank of a tensor is the smallest number of rank-1 tensors that add to \mathcal{T} .

The Norm and Rank of a Tensor



Little difference between matrix and tensor

Consider the tensor \mathcal{T} such that

$$\underline{\underline{\mathcal{T}}} = \underbrace{\mathbf{u} \otimes \mathbf{u} \otimes \mathbf{v}}_{\sim} + \underbrace{\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{u}}_{\sim} + \underbrace{\mathbf{v} \otimes \mathbf{u} \otimes \mathbf{u}}_{\sim}$$

Rank of this tensor seems like 3, it is the limit of these rank-2 tensors \mathcal{T}_n when $n \rightarrow \infty$:

$$\mathcal{T}_n = n \left(\mathbf{u} + \frac{1}{n} \mathbf{v} \right) \otimes \left(\mathbf{u} + \frac{1}{n} \mathbf{v} \right) \otimes \left(\mathbf{u} + \frac{1}{n} \mathbf{v} \right) - n \mathbf{u} \otimes \mathbf{u} \otimes \mathbf{u}$$

Because of the Eckart-Young theorem, the closest approximation to A by a matrix of rank k is fixed. There have been many attempts to decompose tensors, and here are two tensor decompositions. CP and Tucker decomposition.

The CP Decomposition of a Tensor



CP Decomposition

Let $\mathcal{A} \in T_{m,n}$. If there exist a positive integer r , scalars α_j for $j \in \{1, \dots, r\}$, vectors $\mathbf{x}^{(j,i)}$ with $\|\mathbf{x}^{(j,i)}\|^2 = 1$ for $i \in \{1, \dots, m\}$ and $j \in \{1, \dots, r\}$ such that

$$\mathcal{A} = \sum_{j=1}^r \alpha_j \mathbf{x}^{(j,1)} \otimes \dots \otimes \mathbf{x}^{(j,m)}$$

then this summation is said to be a **canonical decomposition/parallel factor decomposition(CANDECOMP/PARAFAC decomposition)** of \mathcal{A} .

It looks similar to SVD but is different. In general, there are many differences between a matrix and a tensor. From the viewpoint of computability, the problem is NP-hard.

The CP Decomposition of a Tensor

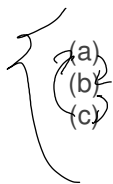


We will consider a 3-order tensor \mathcal{T} .

CP decomposition of \mathcal{T}

$$\mathcal{T} \approx \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c}_1 + \cdots + \mathbf{a}_R \otimes \mathbf{b}_R \otimes \mathbf{c}_R,$$

Here is alternating algorithm by using the three matricized forms T_1, T_2, T_3 .



Minimize $\|T_1 - A(C \odot B)T\|_F^2$

Fix B, C and vary A ,

Fix A, C and vary B ,

Fix A, B and vary C

Matricized Form of a Tensor \mathcal{T}



Suppose A, B, C are the matrices whose columns are the a's and b's and c's. And suppose that the dimension of \mathcal{T} is (I, J, K) . Then each size of matrices is $I \times R, J \times R, K \times R$. We matricize the tensor \mathcal{T} to compute the CP decomposition.

Example

Size = 24

~~| | | | | |
|---|---|---|----|----|
| 1 | 4 | 5 | 10 | 13 |
| 2 | 5 | 7 | 11 | 14 |
| 3 | 6 | 9 | 12 | 15 |~~

Let $I = 3, J = 4, K = 2,$

$$I \times JK = 3 \times 8$$

$$J \times IK = 4 \times 6$$

$$K \times IJ = 2 \times 12$$

$T_1 =$

1	4	7	10	13	16	19	22
2	5	8	11	14	17	20	23
3	6	9	12	15	18	21	24

$T_2 =$

1	2	3	13	14	15
4	5	6	16	17	18
7	8	9	19	20	21
10	11	12	22	23	24

$T_3 =$

1	2	3	4	5	...	9	10	11	12
13	14	15	16	17	...	21	22	23	24

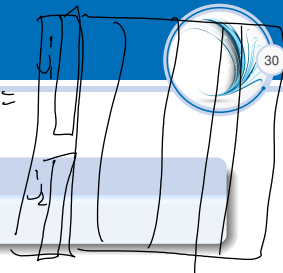
The Khatri-Rao Product $A \odot B$

$$7 + 21 - 2 + 5 = \frac{57}{2}$$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

\times

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



Khatri-Rao

Column j of $A \odot B = (\text{column } j \text{ of } A) \otimes (\text{column } j \text{ of } B)$

Summary

We are slicing T in three directions and placing the slices next to each other in the three matrices T_1, T_2, T_3 . Then we look for three matrices M_1, M_2, M_3 that give us nearly correct equations by ordinary matrix multiplication :

$$\begin{aligned} T_1 &\approx AM_1 & T_2 &\approx BM_2 & T_3 &\approx CM_3 \\ \Rightarrow T_1 &\approx A(C \odot B)^T & T_2 &\approx B(C \odot A)^T & T_3 &\approx C(B \odot A)^T \end{aligned}$$



Recall

Minimize $\|T_1 - A(C \odot B)^T\|_F^2$

- (a) Fix B, C and vary A ,
- (b) Fix A, C and vary B ,
- (c) Fix A, B and vary C

The pseudoinverse of our coefficient matrix $C \odot B$ can be expressed as:

$$(C \odot B)^+ = \underbrace{[(C^T C) \odot (B^T B)]^+}_{\text{}} (C \odot B)^T$$

$$\Rightarrow A = T_1(C \odot B)[(C^T C) \odot (B^T B)]^+$$

Tucker decomposition

Let $\mathcal{A} \in T_{m,n}$. If there exist positive integers r_i for $i \in \{1, \dots, m\}$, scalars $g_{i_1 \dots i_m}$, and vectors $\mathbf{x}^{(j, i_j)}$ with $\|\mathbf{x}^{(j, i_j)}\|_2 = 1$ for $j \in \{1, \dots, m\}$ and $i \in \{1, \dots, n\}$ such that

$$\mathcal{A} = \sum_{i_1=1}^n \cdots \sum_{i_m=1}^n g_{i_1 \dots i_m} \mathbf{x}^{(1, i_1)} \otimes \cdots \otimes \mathbf{x}^{(m, i_m)},$$

then this summation is said to be a Tucker decomposition of \mathcal{A} and the tensor $\mathcal{G} = (g_{i_1 \dots i_m})$ is called the core tensor of \mathcal{A} .

Let $X_j = [\mathbf{x}^{(1, j)} \dots \mathbf{x}^{(r_j, j)}]$ for all j . Then tucker decomposition is related to the k -mode product as follows:

$$\mathcal{A} = \mathcal{G} \times_1 X_1 \times_2 \cdots \times_m X_m.$$

The Tucker Decomposition



Consider the 3-order tensor \mathcal{T} .

Tucker decomposition of \mathcal{T}

$$\mathcal{T} \approx \sum_{p=1}^P \sum_{q=1}^Q \sum_{r=1}^R g_{qpr} \mathbf{a}_p \otimes \mathbf{b}_q \otimes \mathbf{c}_r$$

Similar to the CP decomposition, we can use matricization and tensor product instead of the Khatri-Rao product.

Tucker

$$T_1 \approx \checkmark \underline{AG_1(C \otimes B)^T} \quad T_2 \approx \checkmark \underline{BG_2(C \otimes A)^T} \quad T_3 \approx \checkmark CG_3(B \otimes A)^T$$



Thank you!