

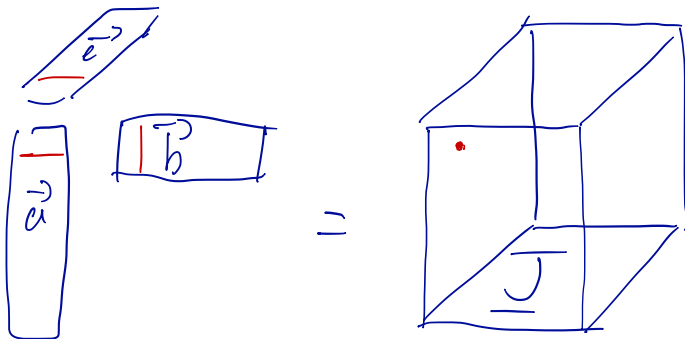
The CP Decomposition and the Rank of a Tensor

Vector outer product of N vectors
 $\vec{a}^{(1)}, \vec{a}^{(2)}, \dots, \vec{a}^{(N)}$ is an N -way tensor

$$\underline{T} = \vec{a}^{(1)} \circ \vec{a}^{(2)} \circ \dots \circ \vec{a}^{(N)}$$

with every element defined as the product of the corresponding elements of the vectors

$$t_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} a_{i_2}^{(2)} \dots a_{i_N}^{(N)}$$



$$t_{ijk} = a_i b_j c_k$$

The CP decomposition

The exact CP decomposition of an N -way tensor $\underline{T} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ has the form

$$\underline{T} = \sum_{r=1}^R \vec{a}_r^{(1)} \circ \vec{a}_r^{(2)} \circ \dots \circ \vec{a}_r^{(N)},$$

where $R \in \mathbb{N}$ and $\vec{a}_r^{(i)} \in \mathbb{R}^{I_i}$ for all $i \in [N]$ and $r \in [R]$.

In the approximate (or fixed-rank) CP decomposition, the size R is given, and we're looking for the least-error decomposition

$$\left\| \underline{T} - \sum_{r=1}^R \vec{a}_r^{(1)} \circ \vec{a}_r^{(2)} \circ \dots \circ \vec{a}_r^{(N)} \right\|.$$

For now, we concentrate on 3-way tensors, and write

$$\underline{T} = \sum_{r=1}^R \vec{a}_r \circ \vec{b}_r \circ \vec{c}_r.$$

Visually, the 3-way CP is

$$T = \vec{a}_1 \vec{b}_1 \vec{c}_1 + \vec{a}_2 \vec{b}_2 \vec{c}_2 + \dots$$

We can gather the vector for each mode in factor matrices. In the 3-way setting, for $T \in \mathbb{R}^{I \times J \times K}$, we have

$$A = [\vec{a}_1 \ \vec{a}_2 \ \dots \ \vec{a}_R], \quad B = [\vec{b}_1 \ \vec{b}_2 \ \dots \ \vec{b}_R], \quad C = [\vec{c}_1 \ \vec{c}_2 \ \dots \ \vec{c}_R]$$

$I \times R$ $J \times R$ $K \times R$

We can express the 3-way CP decomposition using the frontal slices of T and the factor matrices:

$$T_k = A D^{(k)} B^T,$$

where $D^{(k)} = \text{diag}(C(k, :))$, i.e. a diagonal matrix with the k -th row of C on its diagonal.

$$J = A B^T T_k = A D^{(k)} B^T$$

The frontal slice formulation doesn't generalize easily for more than 3 modes. For more generalized representation, we need the Khatri-Rao matrix product: given matrices $A \in \mathbb{R}^{l \times k}$ and $B \in \mathbb{R}^{j \times k}$, their Khatri-Rao product is

$$A \odot B = \begin{pmatrix} a_{11} \vec{b}_1 & a_{12} \vec{b}_2 & \dots & a_{1k} \vec{b}_k \\ a_{21} \vec{b}_1 & a_{22} \vec{b}_2 & \dots & a_{2k} \vec{b}_k \\ \vdots & \vdots & \ddots & \vdots \\ a_{l1} \vec{b}_1 & a_{l2} \vec{b}_2 & \dots & a_{lk} \vec{b}_k \end{pmatrix} \in \mathbb{R}^{l \times k}$$

That is, each column of B is copied l times, and the i -th copy of the k -th column B multiplied by a_{ik} . The Khatri-Rao product can be written more concisely using the

Kronecker matrix product $A \otimes B$
 If $A \in \mathbb{R}^{I \times J}$ and $B \in \mathbb{R}^{K \times L}$, their Kronecker product is

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1J}B \\ a_{21}B & a_{22}B & \dots & a_{2J}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}B & a_{I2}B & \dots & a_{IJ}B \end{pmatrix} \in \mathbb{R}^{(IK) \times (JL)}$$

Notice that in Kronecker product, the matrices can be of arbitrary size, whereas in Khatri-Rao, they must have the same number of columns.

The Khatri-Rao product of $A \in \mathbb{R}^{I \times K}$ and $B \in \mathbb{R}^{J \times K}$ can now be written as

$$A \oslash B = [\vec{a}_1 \otimes \vec{b}_1 \quad \vec{a}_2 \otimes \vec{b}_2 \quad \dots \quad \vec{a}_K \otimes \vec{b}_K],$$

that is, Khatri-Rao is "column-wise Kronecker" product.

If $A \in \mathbb{R}^{I \times R}$, $B \in \mathbb{R}^{J \times R}$, and $C \in \mathbb{R}^{K \times R}$ are the factor matrices of a CP decomposition of tensor $\underline{T} \in \mathbb{R}^{I \times J \times K}$, then

$$T_{(1)} = A(C \circ B)^T$$

$$T_{(2)} = B(C \circ A)^T$$

$$T_{(3)} = C(B \circ A)^T.$$

More generally, if \underline{T} has N modes and factor matrices $A^{(1)}, A^{(2)}, \dots, A^{(N)}$,

$$T_{(n)} = A^{(n)} (A^{(N)} \circ \dots \circ A^{(n+1)} \circ A^{(n-1)} \circ \dots \circ A^{(1)})^T.$$

To gain intuition on the Khatri-Rao formulation, consider the frontal slice formulation of CP:

$$T_k = A D^{(k)} B^T,$$

with $D^{(k)} = \text{diag}(C(:, k))$. The same factor A appears with all frontal slices, so we can just stack them:

$$\begin{bmatrix} T_1 & T_2 & \dots & T_K \end{bmatrix} = A \underbrace{\begin{bmatrix} D^{(1)} B^T & D^{(2)} B^T & \dots & D^{(K)} B^T \end{bmatrix}}_E$$

$$\begin{bmatrix} T_1 & T_2 & \dots & T_k \end{bmatrix} = A \times \underbrace{\begin{bmatrix} D^{(1)} & B^T & D^{(2)} & B^T & \dots & D^{(k)} & B^T \end{bmatrix}}_E$$

The first row of E has the first column of B multiplied by c_{11} followed by the first column of B multiplied by c_{21} , and so on. Hence

$$\vec{e}_{1:} = [c_{11}\vec{b}_1 \quad c_{21}\vec{b}_1 \quad \dots \quad c_{k1}\vec{b}_1] = (\vec{c}_1 \otimes \vec{b}_1)^T$$

Extending this to all rows of E we see that

$$E^T = (C \odot B)^T$$

and hence

$$T_{(1)} = [T_1 \quad T_2 \quad \dots \quad T_k] = A (C \odot B)^T$$

The connections in the other modes can be derived analogously.

One sometimes normalizes the columns of the factor matrices to unit length.

The lengths are then stored in factors

$\lambda_r = \|\vec{a}_r\| \cdot \|\vec{b}_r\| \cdot \|\vec{c}_r\|$, collected in a vector $\vec{\lambda} \in \mathbb{R}^R$, or in a matrix $\Lambda = \text{diag}(\vec{\lambda}) \in \mathbb{R}^{R \times R}$.

Then $\hat{T}_{(1)} = A \Lambda (C \circ B)^T$ etc. and

$$\underline{T} = \sum_{r=1}^R \lambda_r \vec{a}_r \circ \vec{b}_r \circ \vec{c}_r.$$

A common notation for the CP decomposition is to write

$$\underline{T} = [A, B, C] = \sum_{r=1}^R \vec{a}_r \circ \vec{b}_r \circ \vec{c}_r,$$

or with the scaling

$$\underline{T} = [\underline{\lambda}; A, B, C] = \sum_{r=1}^R \lambda_r \vec{a}_r \circ \vec{b}_r \circ \vec{c}_r.$$

ALS algorithm for CP

The formulations

$$T_{(1)} = A(C \circ B)^T \text{ etc}$$

provide a way to solve the (approximate) CP decomposition. When C and B are fixed $(C \circ B)$ is a fixed matrix, call it D , and the problem becomes: "Given matrices $T_{(1)}$ and D , find matrix A that minimizes $\|T_{(1)} - AD^T\|_F$ ". This can be solved using the SVD and pseudo-inverse as $A = T_{(1)}(D^T)^+$, where $(\cdot)^+$ is the Moore-Penrose pseudo-inverse. This leads to the following algorithm:

sample random B and C

repeat

$$\text{let } A \leftarrow T_{(1)}((C \circ B)^T)^+$$

$$\text{let } B \leftarrow T_{(2)}((C \circ A)^T)^+$$

$$\text{let } C \leftarrow T_{(3)}((B \circ A)^T)^+$$

until convergence

The ALS algorithm requires us to compute the pseudo-inverses of $(C \odot B)^T$, $(C \odot A)^T$, and $(B \odot A)^T$, which are R -by- Jk , R -by- Ik , and R -by- IJ matrices, respectively. This is an expensive operation, but if these matrices have a full row rank — which is likely, as often $R \ll \min\{IJ, Ik, Jk\}$ — then we can use the following equality

$$(A \odot B)^+ = ((A^T A) * (B^T B))^+ (A \odot B)^T, \quad (*)$$

where $X * Y$ is the Hadamard matrix product (or elementwise product) between $X \in \mathbb{R}^{I \times J}$ and $Y \in \mathbb{R}^{I \times J}$

$$X * Y = \begin{pmatrix} x_{11} y_{11} & x_{12} y_{12} & \dots & x_{1j} y_{1j} \\ x_{21} y_{21} & x_{22} y_{22} & \dots & x_{2j} y_{2j} \\ \vdots & \vdots & \ddots & \vdots \\ x_{i1} y_{i1} & x_{i2} y_{i2} & \dots & x_{ij} y_{ij} \end{pmatrix} \in \mathbb{R}^{I \times J}$$

The proof of identity $(*)$ is left as a homework, but it involves the following identity that is also occasionally useful on itself:
 For $X \in \mathbb{R}^{l \times k}$ and $Y \in \mathbb{R}^{j \times k}$, we have

$$(X \circ Y)^T (X \circ Y) = X^T X * Y^T Y.$$

Proof: Let $X' = X^T X$ and notice that $x'_{ik} = \langle \vec{x}_i, \vec{x}_k \rangle$. Similarly, for $Y' = Y^T Y$, we have $y'_{jk} = \langle \vec{y}_j, \vec{y}_k \rangle$. Now, let

$$\begin{aligned} Z &= (X \circ Y)^T (X \circ Y) \\ &= [\vec{x}_1 \otimes \vec{y}_1, \dots, \vec{x}_k \otimes \vec{y}_k]^T [\vec{x}_1 \otimes \vec{y}_1, \dots, \vec{x}_k \otimes \vec{y}_k], \end{aligned}$$

and consider a single element z_{kl} :

$$\begin{aligned} z_{kl} &= \langle \vec{x}_k \otimes \vec{y}_k, \vec{x}_l \otimes \vec{y}_l \rangle = \sum_{i=1}^l \langle x_{ik} \vec{y}_k, x_{il} \vec{y}_l \rangle \\ &= \sum_{i=1}^l x_{ik} x_{il} \langle \vec{y}_k, \vec{y}_l \rangle \\ &= \langle \vec{x}_k, \vec{x}_l \rangle \langle \vec{y}_k, \vec{y}_l \rangle = x'_{kl} y'_{kl}, \end{aligned}$$

and hence $Z = X' * Y' = X^T X * Y^T Y$. \square

With equality (*) we can write

$$A = T_{(1)} \underbrace{\left[(C \circ B)^T \right]}_{R \times JK}^+$$

as

$$A = T_{(1)} (C \circ B) \underbrace{\left[C^T C * B^T B \right]}_{R \times R}^+,$$

and we only have to take the pseudo-inverse from a much smaller matrix. This formulation can, however, cause issues with the numerical stability.

ALS is not the only possibility. We can instead use, for instance, gradient-based methods: each row $\vec{a}_{i:}$ can be updated based on the gradient

$$\vec{a}_{i:} \leftarrow \vec{a}_{i:} - \delta \frac{\partial}{\partial \vec{a}_{i:}} \sum_{j=1}^{JK} \left(T_{(1)}(i, j) - (\vec{a}_{i:} (C \circ B)^T)_j \right)^2.$$

ALS is the most commonly used approach, though.