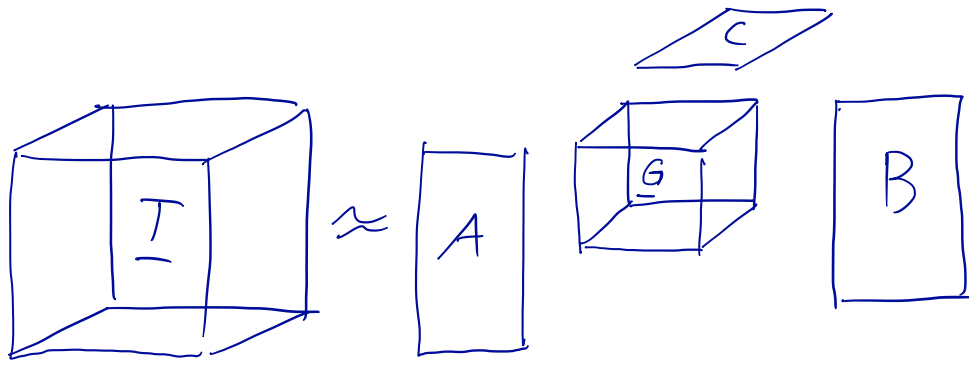# The Tucker Decompositions

Let $\underline{T} \in \mathbb{R}^{I \times J \times K}$. Its Tucker3 — or just Tucker — decomposition has a _core tensor_ $\underline{G} \in \mathbb{R}^{P \times Q \times R}$ and thee factor matrices $A \in \mathbb{R}^{I \times P}$, $B \in \mathbb{R}^{J \times Q}$, and $C \in \mathbb{R}^{K \times R}$, and it is defined as

$$\underline{T} \approx \underline{G} \times_1 A \times_2 B \times_3 C$$

$$= \sum_{p=1}^{P} \sum_{q=1}^{Q} \sum_{r=1}^{R} g_{pqr} \, \vec{a}_p \circ \vec{b}_q \circ \vec{c}_r$$

$$= [\![ \underline{G}; A, B, C ]\!] .$$

Hence, elementwise Tucker is

$$t_{ijk} \approx \sum_{p=1}^{P} \sum_{q=1}^{Q} \sum_{r=1}^{R} g_{pqr} \, a_{ip} \, b_{jq} \, c_{kr} .$$

Tucker decomposition is usually applied only to 3-way tensors, but N-way version is straightforward to define:

$$\underline{T} \approx \underline{G} \times_1 A^{(1)} \times_2 A^{(2)} \times_3 \cdots \times_N A^{(N)}$$

$$t_{i_1 i_2 \cdots i_N} \approx \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \cdots \sum_{r_N=1}^{R_N} g_{r_1 r_2 \cdots r_N} \, a_{i_1 r_1}^{(1)} \, a_{i_2 r_2}^{(2)} \cdots a_{i_N r_N}^{(N)}$$

The core tensor $\underline{G}$ can be considered as a compressed version of the original tensor $\underline{T}$ if $P < I$, $Q < J$, and $R < K$. If $P = Q = R$ and $\underline{G}$ is (hyper-) diagonal, then Tucker3 reduces to the CP decomposition. In particular, with hyper-diagonal $\underline{G}$ where all diagonal entries are 1, we have that $[\![ \underline{G}; A, B, C ]\!] = [\![ A, B, C ]\!]$.

The Tucker decomposition can be expressed in a matricized form:

$$T_{(1)} \approx A G_{(1)} (C \otimes B)^T$$
$$T_{(2)} \approx B G_{(2)} (C \otimes A)^T$$
$$T_{(3)} \approx C G_{(3)} (B \otimes A)^T.$$

To gain some intuition, let's consider the first frontal slice of $[\![G; A, B, C]\!]$.

$$[\![\underline{G}; A, B, C]\!]_1 = A D^{(1)} B^T, \text{ with } D^{(1)} = \sum_{r=1}^{R} c_{1r} G_r.$$

Now we have

$$T_{(1)} = A [D^{(1)} B^T \; D^{(2)} B^T \; \cdots \; D^{(k)} B^T]$$
$$= A [(\sum_r c_{1r} G_r) B^T \; \cdots \; (\sum_r c_{kr} G_r) B^T]$$
$$= A [G_1 \; G_2 \; \cdots \; G_k] \begin{pmatrix} c_{11} B^T & c_{21} B^T & \cdots & c_{k1} B^T \\ c_{12} B^T & c_{22} B^T & \cdots & c_{k2} B^T \\ \vdots & \vdots & \ddots & \vdots \\ c_{1R} B^T & c_{2R} B^T & \cdots & c_{kR} B^T \end{pmatrix}$$
$$= A G_{(1)} (C \otimes B)^T$$
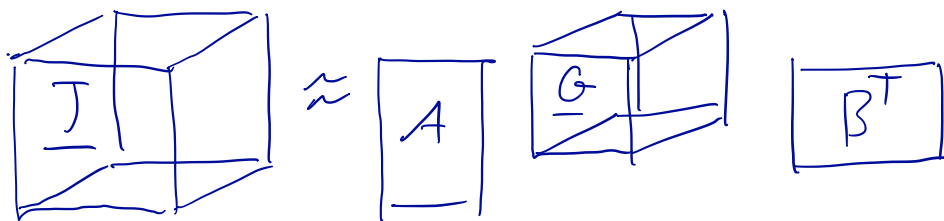
In N-way case, the mode-$n$ matricized version becomes

$$T_{(n)} \approx A^{(n)} G_{(n)} \left( A^{(N)} \otimes \cdots \otimes A^{(n+1)} \otimes A^{(n-1)} \otimes \cdots \otimes A^{(1)} \right)^T$$

## Tucker1 and Tucker2

The Tucker2 decomposition leaves one factor matrix as an identity matrix, e.g.

$$\underline{T} \approx \underline{G} \times_1 A \times_2 B = [\![ \underline{G}; A, B, I ]\!],$$

where $\underline{T} \in \mathbb{R}^{(x) \times k}$ and $\underline{G} \in \mathbb{R}^{P \times Q \times k}$



The Tucker1 decomposition leaves two factor matrices as identity, for instance

$$\underline{T} \approx \underline{G} \times_1 A = [\![ \underline{G}; A, I, I ]\!].$$

This is equivalent to standard least-squares matrix factorization.

# The n-rank

Let $T \in \mathbb{R}^{I_1 \times I_2 \times \cdots I_N}$. The n-rank of $T$, $\mathrm{rank}_n(T)$, is the column rank of $T_{(n)}$, i.e., the number of linearly independent columns in $T_{(n)}$. If we set $R_n = \mathrm{rank}_n(T)$, then $T$ is rank-$(R_1, R_2, \ldots, R_N)$, though note that this definition is not compatible with the usual tensor rank.

(Clearly, $R_n \leq I_n$ for all $n \in [N]$. Finding a Tucker decomposition of size $(R_1, R_2, \ldots, R_N)$ is easy if $\mathrm{rank}_n(T) = R_n$ for all $n$.

# Computing Tucker

From now on, we enforce that the columns of the factor matrices are mutually orthogonal, that is,
$$A^T A = I_P, \quad B^T B = I_Q, \quad \text{and} \quad C^T C = I_R.$$

# Higher-order SVD (HOSVD)

HOSVD is a simple method to calculate the Tucker decomposition using SVD:

$$\text{for } n = 1, \dots, N \text{ do}$$
$$(U, \Sigma, V) \leftarrow \text{SVD}(T_{(n)})$$
$$A^{(n)} \leftarrow U(:, 1:R_n)$$
$$\text{end}$$
$$G \leftarrow T \times_1 A^{(1)T} \times_2 A^{(2)T} \times_3 \cdots \times_N A^{(N)T}$$

In the last row above, we use the fact that the factor matrices are column-orthogonal:

$$T = G \times_1 A \iff T_{(1)} = A G_{(1)} \iff A^{-1} T_{(1)} = G_{(1)}$$
$$\overset{*}{\iff} A^T T_{(1)} = G_{(1)} \iff T \times_1 A^T = G_{(1)},$$

where $*$ is due to the column orthogonality.

HOSVD is not very good for finding a low-error decomposition, as it doesn't take into account the higher-order structure in the tensor. It is, however, useful for finding an initial solution for other algorithms.

# Higher-order Orthogonal Iteration (HOOI)

We know that
$$G = T \times_1 A^T \times_2 B^T \times_3 C^T \qquad (x)$$
is optimal for column-orthogonal $A, B,$ and $C$.
Let us re-write the objective:

$$\| T - [\![ G; A, B, C ]\!] \|^2 = \| T \|^2 - 2 \langle T, [\![ G; A, B, C ]\!] \rangle$$
$$+ \| [\![ G; A, B, C ]\!] \|^2$$
$$\overset{*}{=} \| T \|^2 - 2 \langle T, [\![ G; A, B, C ]\!] \rangle + \| G \|^2$$
$$\overset{**}{=} \| T \|^2 - 2 \langle T \times_1 A^T \times_2 B^T \times_3 C^T, G \rangle + \| G \|^2$$
$$\overset{***}{=} \| T \|^2 - 2 \langle G, G \rangle + \| G \|^2$$
$$\overset{****}{=} \| T \|^2 - \| G \|^2$$
$$= \| T \|^2 - \| T \times_1 A^T \times_2 B^T \times_3 C^T \|^2 .$$

* We use the column orthogonality

** We multiply the factors out from $[\![ G; A, B, C ]\!]$.

*** We use $(x)$

**** We use $\langle G, G \rangle = \| G \|^2$

As $\| T \|$ is constant, we learn that

$$\| T - [\![ G; A, B, C ]\!] \|^2 \propto - \| T \times_1 A^T \times_2 B^T \times_3 C^T \|^2$$

Hence, we want to _maximize_

$$\max_{A,B,C} \| \underline{I} \times_1 A^T \times_2 B^T \times_3 C^T \|,$$

or equivalently

$$\max_{A,B,C} \| A^T T_{(1)} (C \otimes B) \|$$

$$\max_{A,B,C} \| B^T T_{(2)} (C \otimes A) \|$$

$$\max_{A,B,C} \| C^T T_{(3)} (B \otimes A) \|.$$

Essentially, $\underline{G}$ is fully determined by the factor matrices, and we do not have to take it into account.

To compute HOOI, we again use SVD:

Initialize $A^{(1)}, A^{(2)}, \ldots$ using HOSVD
repeat
    for $n = 1, \ldots, N$ do
        $\underline{Y} \leftarrow \underline{I} \times_1 A^{(1)T} \times_2 \cdots \times_N A^{(N)T}$
        $U, \Sigma, V \leftarrow \text{SVD}(Y_{(n)})$
        $A^{(n)} \leftarrow U(:, 1:R_n)$
    end
until convergence
$\underline{G} \leftarrow \underline{I} \times_1 A^{(1)T} \times_2 \cdots \times_N A^{(N)T}$

# Naive ALS for Tucker

We can also use the standard ALS approach to solve Tucker decomposition. This does not, generally, yield te orthogonal factor matrices, though. In some cases, this is what we want. To update the factor matrices, we have:

$$A \leftarrow T_{(1)} \left( G_{(1)} \left( (C \otimes B)^T \right) \right)^+$$
$$B \leftarrow T_{(2)} \left( G_{(2)} \left( (C \otimes A)^T \right) \right)^+$$
$$C \leftarrow T_{(3)} \left( G_{(3)} \left( (B \otimes A)^T \right) \right)^+$$

To update the core, we can use a vectorized format of Tucker and solve

$$G = \arg\min_{G} \| \mathrm{vec}(T) - (C \otimes B \otimes A) \mathrm{vec}(G) \|,$$

which is just a standard least-squares problem.

The ALS algorithm is rarely used due to the large pseudo-inverse.

It can be useful for computing the Tucker2 decomposition, though. For Tucker 2, we replace $C$ with the identity matrix, and solve

$$A \leftarrow T_{(1)} \left( G_{(1)} (I \otimes B)^T \right)^+$$
$$B \leftarrow T_{(2)} \left( G_{(2)} (I \otimes A)^T \right)^+,$$

where we can use the fact that

$$(I \otimes B)^T = \begin{pmatrix} B^T & & & \\ & B^T & & \\ & & \ddots & \\ & & & B^T \end{pmatrix}$$

and that

$$G_{(1)} (I \otimes B)^T = [G_1 \ G_2 \cdots G_k] (I \otimes B)^T = \begin{pmatrix} G_1 B^T & & & \\ & G_2 B^T & & \\ & & \ddots & \\ & & & G_k B^T \end{pmatrix}$$

The core can be updated separately for each frontal slice:

$$G_k \leftarrow B^+ T_k A^+.$$

# Nonnegative Tucker3

Similarly to NCP, we can consider the nonnegative variant of the Tucker3 decomposition. We use the unfolded versions of Tucker

$$T_{(1)} \approx A G_{(1)} (C \otimes B)^T$$
$$T_{(2)} \approx B G_{(2)} (C \otimes A)^T$$
$$T_{(3)} \approx ( G_{(3)} (B \otimes A)^T,$$

and update the factors using multiplicative rules similar to those in NCP.

To initialize $A, B, C,$ and $G$ we can run HOSVD on $T$ and truncate the negative values in the result to zero. We can also use random initialization, but this can yield very bad initial solutions.

The nonnegative Tucker3 algorithm is

Input: $\underline{T} \in \mathbb{R}_{\geq 0}^{I_1 \times I_2 \times \cdots \times I_N}$, $(R_1, R_2, \ldots, R_N)$

$(\underline{G}, A^{(1)}, \ldots, A^{(N)}) \leftarrow \text{HOSVD}(\underline{T}, R_1, R_2 \cdots, R_N)$

$\underline{G} \leftarrow [\underline{G}]_+$ ; $A^{(n)} \leftarrow [A^{(n)}]_+$ for all $n \in [N]$

repeat

$\quad \underline{Q} \leftarrow \underline{G} \times_1 A^{(1)} \times_2 \cdots \times_N A^{(N)}$

$\quad$ for $n = 1 \ldots N$ do

$\quad\quad A^{(n)} \leftarrow A^{(n)} * \left( T_{(n)} A^{\otimes_{-n}} G_n^{\top} \right) \oslash \left( Q_{(n)} A^{\otimes_{-n}} G_{(n)}^{\top} \right)$

$\quad\quad \vec{a}_i^{(n)} \leftarrow \vec{a}_i^{(n)} / \| \vec{a}_i^{(n)} \|$

$\quad$ end

$\quad \underline{G} \leftarrow \underline{G} * \left( \underline{T} \times_1 A^{(1)} \times_2 \cdots \times_N A^{(N)} \right) \oslash \left( Q_{(n)} \times_1 A^{(1)} \times_2 \cdots \times_N A^{(N)} \right)$

until convergence

Here, $A^{\otimes_{-n}} = A^{(N)} \otimes \cdots \otimes A^{(n+1)} \otimes A^{(n-1)} \otimes \cdots \times A^{(1)}$.

# Applications of Tucker

## Tensor Faces

We can use Tucker decomposition to separate the effects of different illuminations, expressions, poses, etc from pictures, provided that we have training data with all variations for all subjects.

In Tensor Faces, we have photos of subjects in different views (front, left, right,...), under different illuminations and having different expressions. This gives a people-by-views-by-illuminations-by-expressions-by-pixels tensor. See the next page for some examples (image from Vasilescu & Terzopoulos: Multilinear analysis of image ensembles: Tensor Faces, ECCV'02).

Applying the Tucker decomposition with no reduction in the dimensions yields

$$\underline{G} \times_1 A_{people} \times_2 A_{views} \times_3 A_{illums} \times_4 A_{express} \times_5 A_{pixels}$$

The $A_{(n)}$ factor matrices encode the variations in these modes, while the core tensor governs the interactions.

This model generalizes Eigenfaces: we can write

$$T_{(pixels)} = A_{pixels} G_{(pixels)} \left( A_{express} \otimes A_{illums} \otimes A_{views} \otimes A_{people} \right)^T$$

to get a standard Eigenfaces setting. Multiplying $\underline{G} \times_5 A_{pixels}$ gives us a tensor that shows the primary variations along the modes. Some examples are in the next page (Vasilescu & Terzopoulos, 2002). On the other hand, if we multiply

$$\underline{G} \times_2 A_{views} \times_5 A_{pixels}$$

We get the variations with the different viewpoints, as depicted in the picture in the previous page (Vasilescu & Terzopoulos, 2002).

## Finding facts in openIR

In open information retrieval our goal is to extract structured knowledge from unstructured data using unsupervised methods. Typically, we want to extract subject-predicate-object $(s, p, o)$ triples with disambiguated entities and relations. To do that, we can first run standard natural language parsers to obtain noun phrase-verbal phrase-noun phrase $(np, vp, np)$ triples, e.g

Donald J. Trump, is, POTUS
Donald Trump, is the president of, USA
The Donald, is the prez of, 'Murica
Donald Trump, is the son of, POTUS

These triples encode two $(S, P, O)$ triples

   donald_j_trump, isPresidentOf, USA
   donald_j_trump_jr, isSonOf, donald_j_trump

To extract these, we need to handle synonyms (president vs. prez) and homonyms (Donald Trump [Jr]), among others. We can model this with mappings from noun phrases to entities and from verbal phrases to relations. Considering subjects and objects separately, we have $A: np \rightarrow s$, $B: np \rightarrow o$, and $C: vp \rightarrow p$. We model these as matrices. If they are column orthogonal, and if our original $(np, vp, np)$ triples are stored in a tensor $\underline{T}$, we can get the true $(s, p, o)$ triples as

$$\underline{T} \times_1 A^T \times_2 B^T \times_3 C^T.$$

That is, if we do Tucker decompo-
sition to the tensor containing the
surface triples, we find the latent
entities and relations (though they
might not contain any "real"
entities or relations).

A practical problem with this approach
is that the core tensor has to be very
big, making it hard to work with.
The core should also be sparse, which
is not the case if we obtain it with
the multiplication.