

# Tensors in Data Analysis

15 May 2014

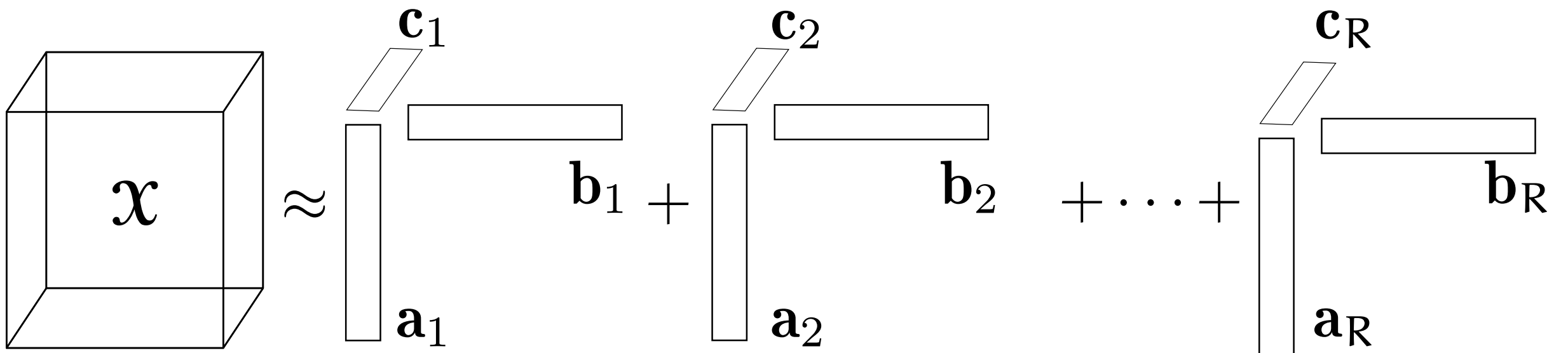


**mpi** max planck institut  
informatik

# Tensors in Data Analysis

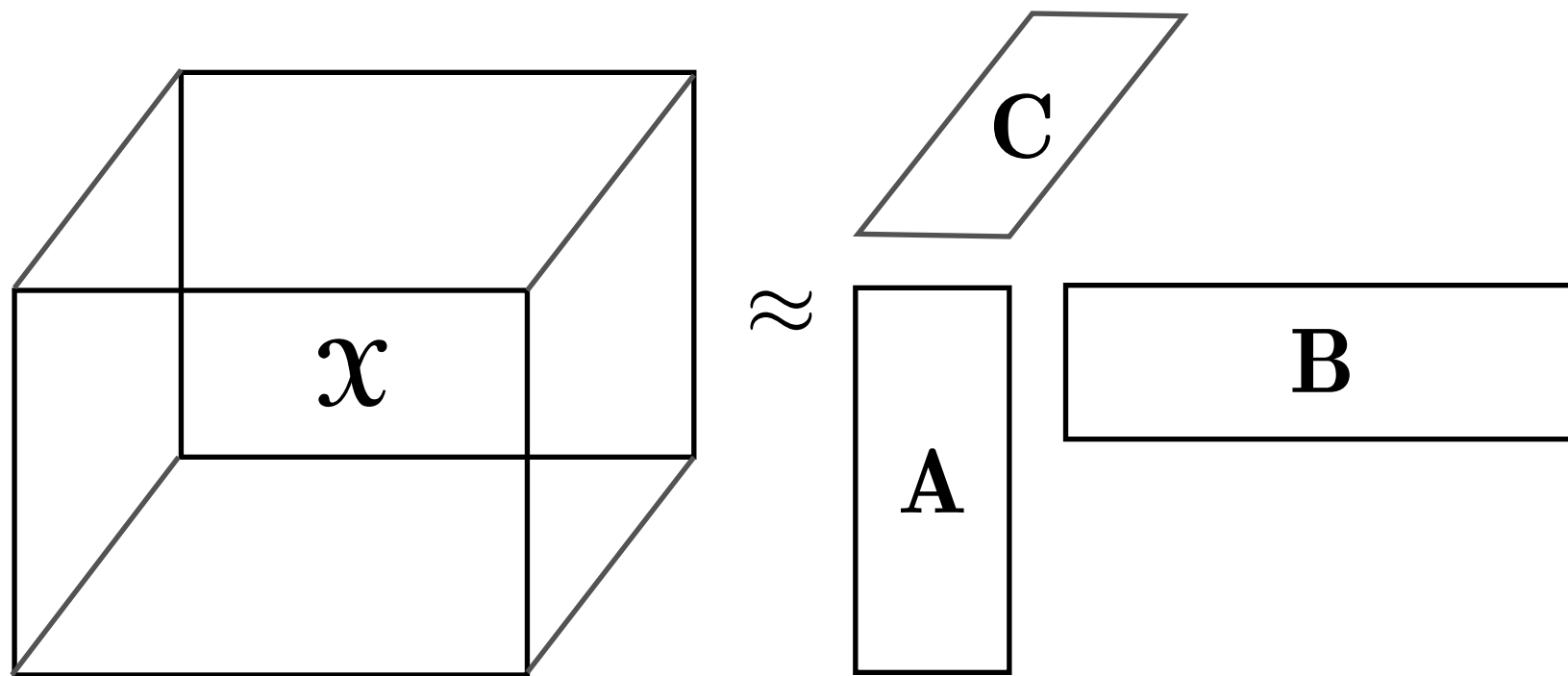
1. CP and INDSCAL and some applications
2. The Tucker tensor decompositions
3. HOSVD, RESCAL, and DEDICOM
4. The non-negative variants

# CP Recap (Rank-1 View)



$$x_{ijk} \approx \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$$

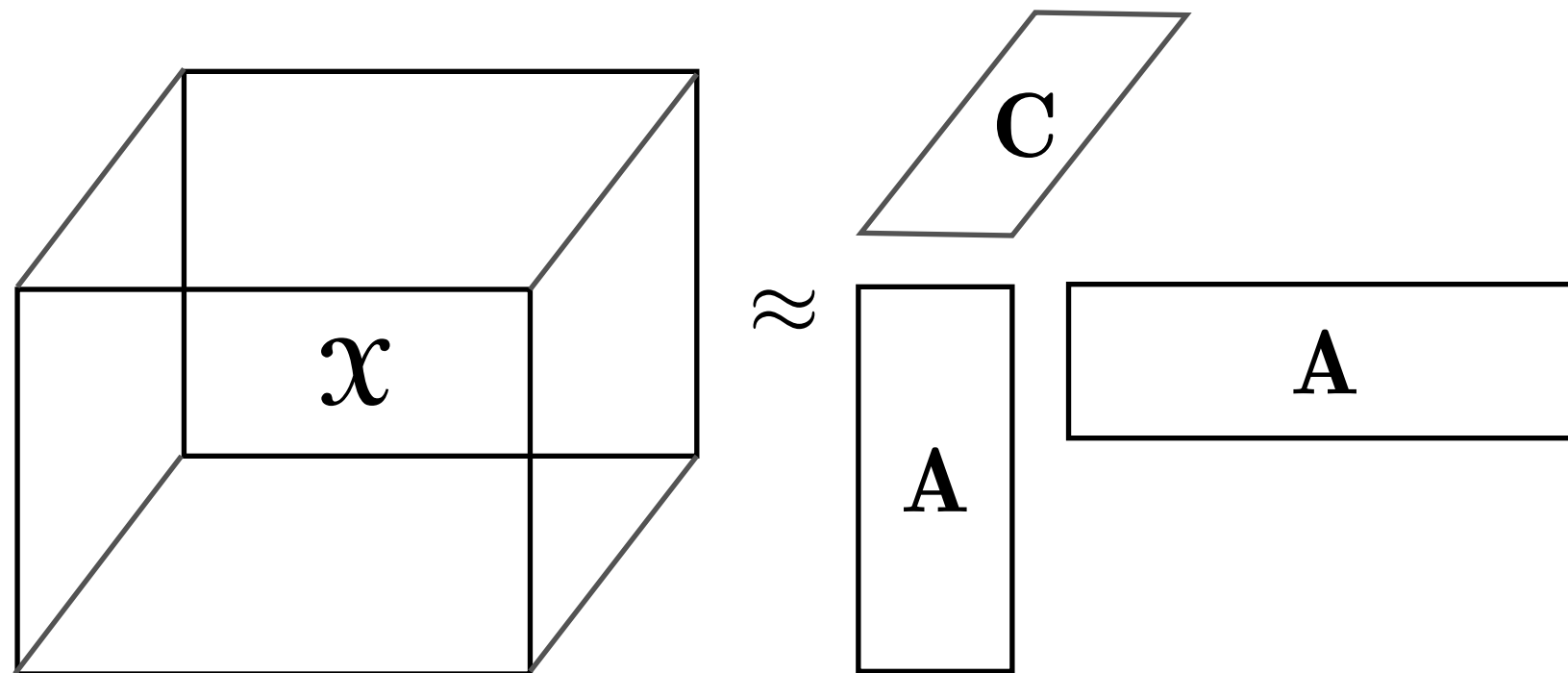
# CP Recap (Matrix View)



$$x_{ijk} \approx \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$$

# The INDSCAL Decomposition

- The **INDSCAL decomposition** decomposes a 3-way tensor  $\mathcal{X}$  into **two** factor matrices  **$\mathbf{A}$**  and  **$\mathbf{C}$**



$$x_{ijk} \approx \sum_{r=1}^R a_{ir} a_{jr} c_{kr}$$

# More on INDSCAL

- First two modes of  $X$  are expected to be symmetric
  - Not mandatory, but must have same dimensions
- Commonly computed by solving CP and hoping **A** and **B** merge
  - End by forcing **A** and **B** the same and update **C**

# Why INDSCAL

- INDSCAL keeps the symmetry of the modes
- Stands for Individual Differences in Scaling
  - Assume  $K$  subjects ranked the similarity of  $N$  objects
  - Assume each subject is influenced by the same factors, but with different weights
  - **$A$**  contains the factors,  **$C$**  gives the weights

# INDSCAL Example

- Carroll and Chang (1970) proposed to use INDSCAL and CANDECOMP (CP) to analyse psychological data
  - PCA has long history in psychology
- Example: 20 subjects rate the similarity of countries
  - Multi-way data



# Countries in Carroll & Chan (1970) [1]

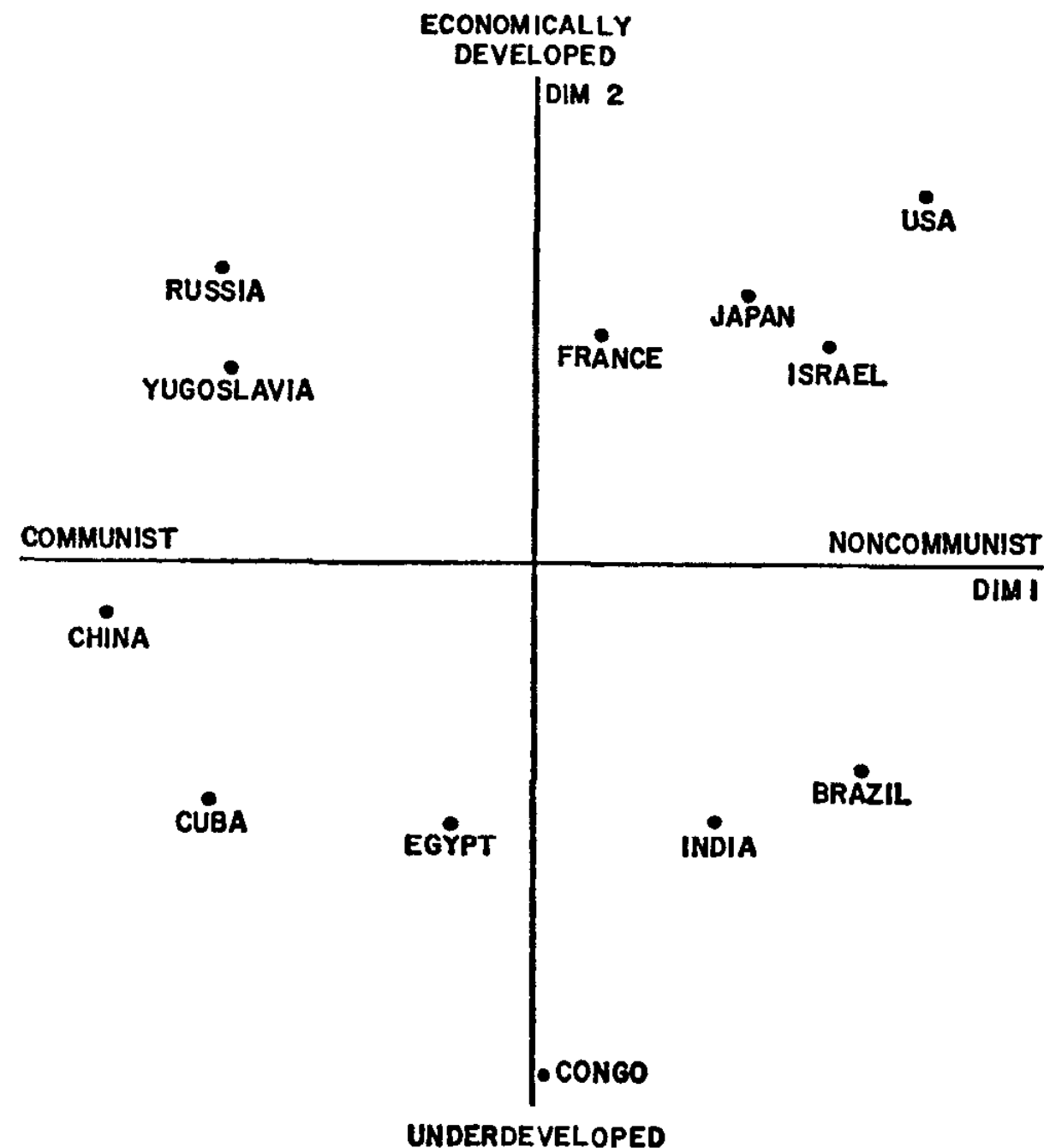


FIGURE 11

# Countries in Carroll & Chan (1970) [2]

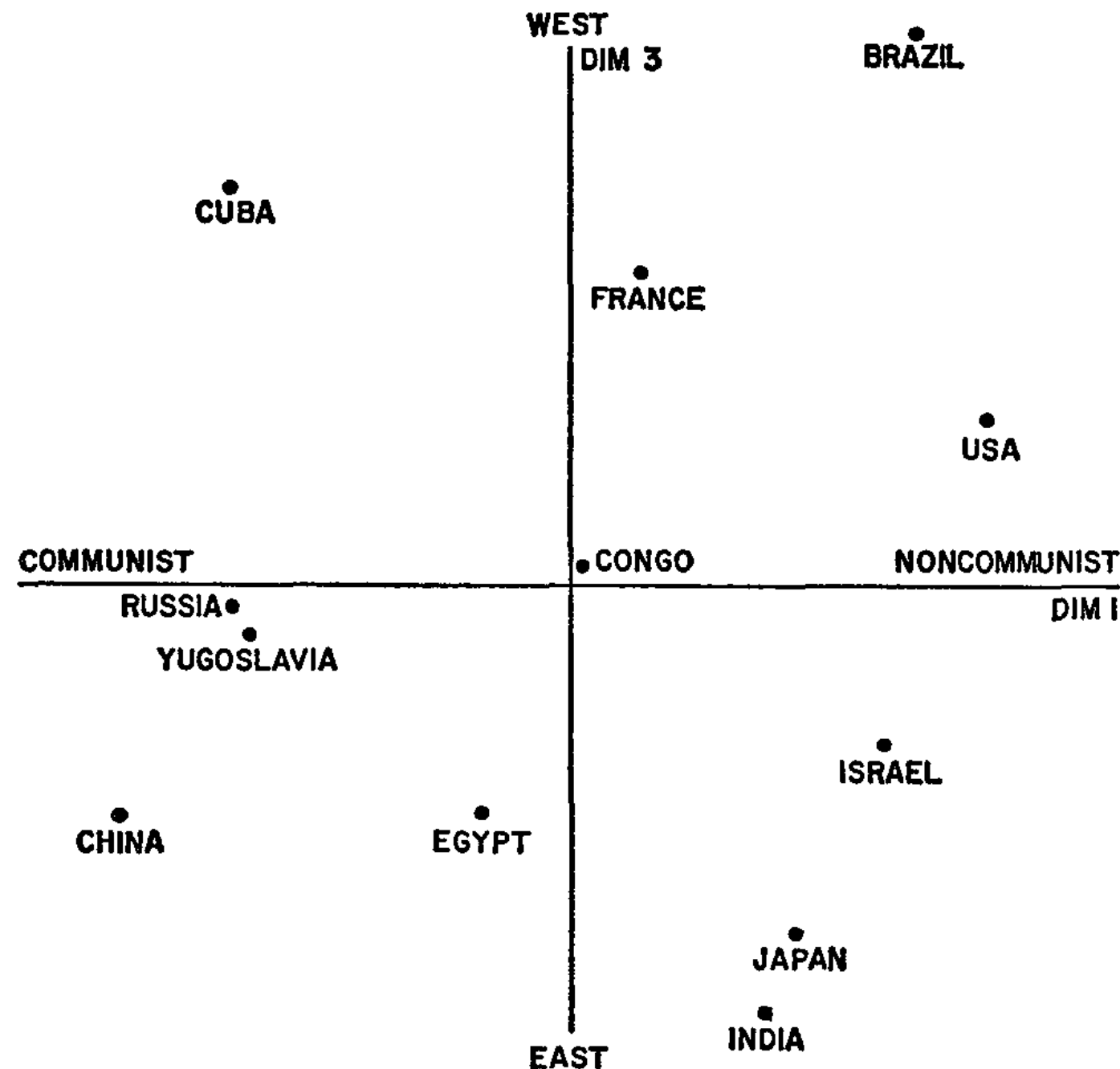


FIGURE 12

# Countries in Carroll & Chan (1970) [3]

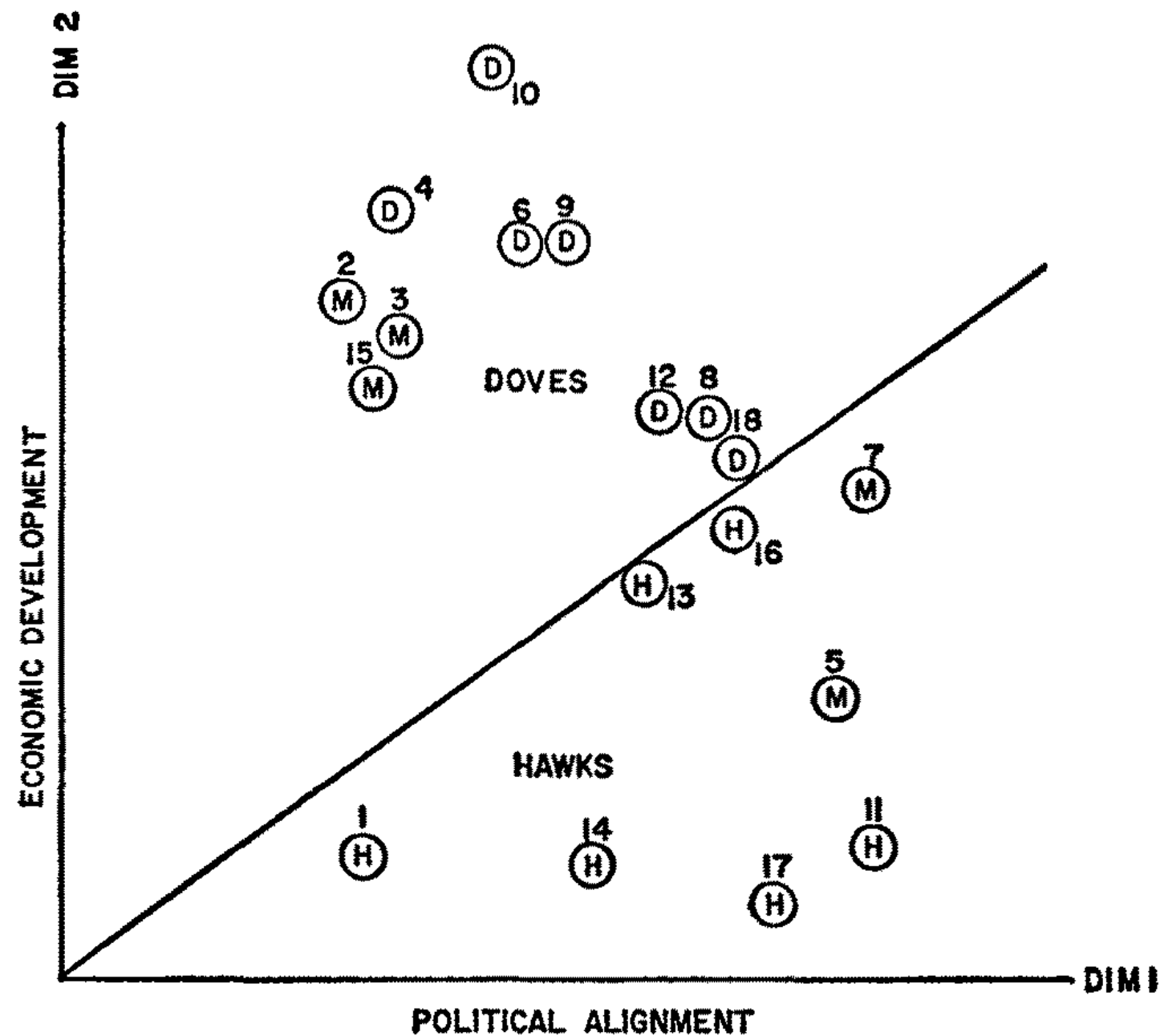


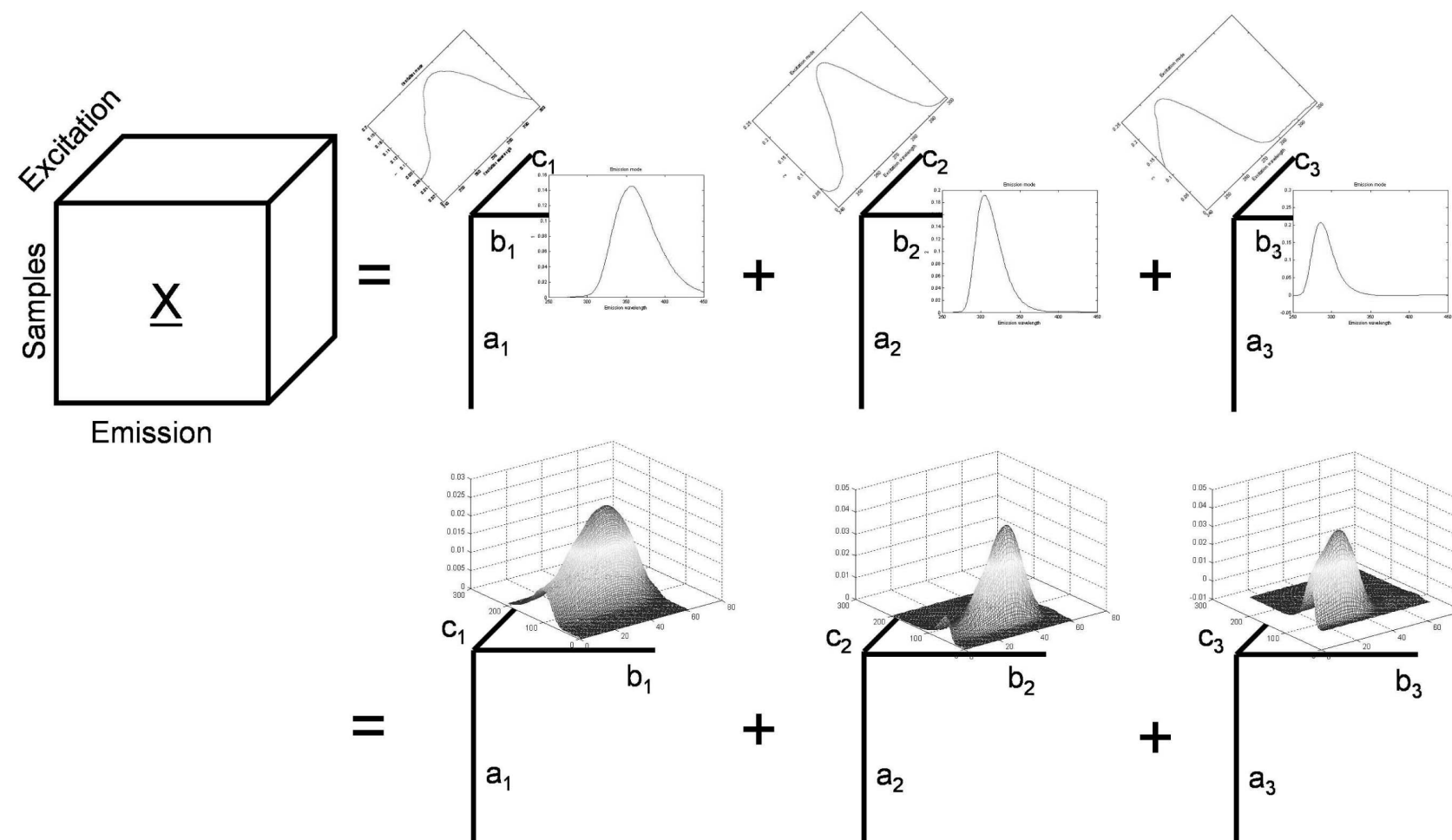
FIGURE 13

# Fluorescence Excitation– Emission Analysis

- Fluorescence spectroscopy analyses (typically) organic compounds
  - A beam of (UV) light excites electrons in molecules
  - The excited electrons release a photon, which is measured
  - A **fluorescence landscape** of a compound is a rank-1 matrix that maps the exciter's wavelength to the emitted photon's wavelength
    - Lets us to identify the compounds

# CP for Fluorescence Analysis

- Samples-by-excitation wavelengths-by-emission wavelengths tensor  $\mathcal{X}$ 
  - Matrix  $\mathbf{b}_i \mathbf{c}_i^T$  is the landscape for the  $i$ th component
  - Vector  $\mathbf{a}_i$  gives the weights of landscapes in each sample



# TOPHITS for IR

- Three-way pages-by-pages-by-anchor text tensor  $\mathcal{T}$ 
  - Element  $t_{ijk} = \max\{1 + \log(x_{ijk}), 0\}$  where  $x_{ijk}$  is the number of times page  $i$  links to page  $j$  using term  $k$
- The CP decomposition of  $\mathcal{T}$  behaves akin to HITS
  - Each rank-1 component is one topic
  - **A** and **B** give the authority and hub scores, **C** gives the weights for terms

# The Tucker Decompositions

- The CP decomposition requires the factors to have the same number of columns
- In Tucker decompositions, different number of columns can be mixed using a **core tensor**
  - This enables very different looking decompositions

# Tensor–Vector Multiplication

- Vectors can be multiplied with tensors along specific modes
  - For  $n$ -th mode multiplication, the tensor's dimensionality in mode  $n$  must agree with the vector's dimensions
- The  $n$ -mode vector product is denoted  $\mathcal{X} \bar{\mathbf{x}}_n \mathbf{v}$ 
  - The result is of order  $N-1$
  - $(\mathcal{X} \bar{\mathbf{x}}_n \mathbf{v})_{i_1 \cdots i_{n-1} i_{n+1} \cdots i_N} = \sum_{i_n=1}^{I_n} \mathcal{X}_{i_1 i_2 \cdots i_N} \mathbf{v}_{i_n}$ 
    - Inner product between mode- $n$  fibres and vector  $\mathbf{v}$



# Tensor-Vector Multiplication Example

Given tensor  $\mathcal{T}$  and vector  $\mathbf{v}$ ,

$$\mathbf{T}_1 = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \quad \mathbf{T}_2 = \begin{pmatrix} 5 & 7 \\ 6 & 8 \end{pmatrix} \quad \mathbf{v} = \begin{pmatrix} 2 & 1 \end{pmatrix}$$

Computing  $\mathcal{Y} = \mathcal{T} \bar{\times}_3 \mathbf{v}$  gives

$$\mathcal{Y} = \begin{pmatrix} 7 & 13 \\ 10 & 16 \end{pmatrix}$$

# Tensor–Matrix Multiplication

- Let  $\mathcal{X}$  be an  $N$ -way tensor of size  $I_1 \times I_2 \times \dots \times I_N$ , and let  $\mathbf{U}$  be a matrix of size  $J \times I_n$
- The  $n$ -mode matrix product of  $\mathcal{X}$  with  $\mathbf{U}$ ,  $\mathcal{X} \times_n \mathbf{U}$  is of size  $I_1 \times I_2 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N$
- $(\mathcal{X} \times_n \mathbf{U})_{i_1 \dots i_{n-1} j i_{n+1} \dots i_N} = \sum_{i_n=1}^{I_n} \mathcal{X}_{i_1 i_2 \dots i_N} u_{j i_n}$ 
  - Each mode- $n$  fibre is multiplied by the matrix  $\mathbf{U}$
- In terms of unfold tensors:

$$\mathcal{Y} = \mathcal{X} \times_n \mathbf{U} \iff \mathbf{Y}_{(n)} = \mathbf{U} \mathbf{X}_{(n)}$$

# Tensor-Matrix Multiplication Example

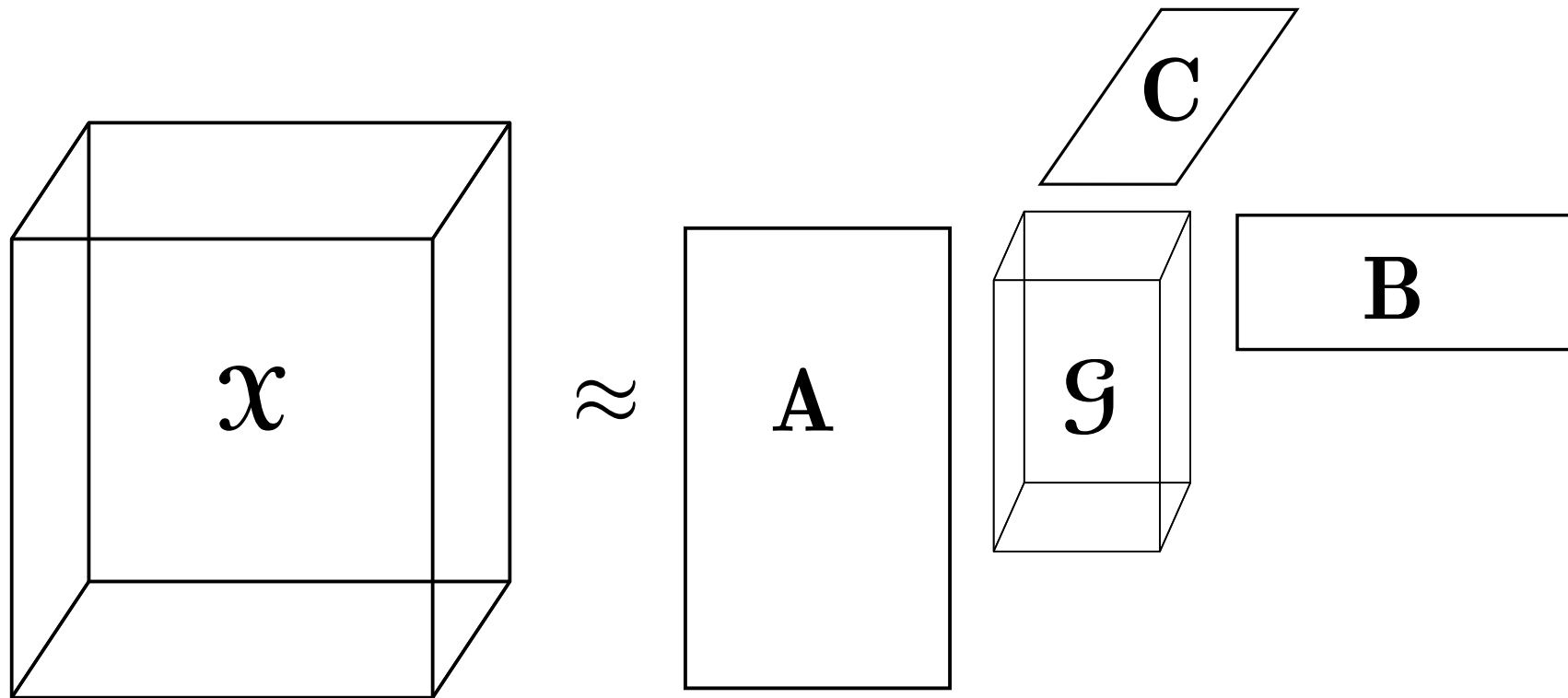
Given tensor  $\mathcal{T}$  and matrix  $\mathbf{M}$ ,

$$\mathbf{T}_1 = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \quad \mathbf{T}_2 = \begin{pmatrix} 5 & 7 \\ 6 & 8 \end{pmatrix} \quad \mathbf{M} = \begin{pmatrix} 10 & 0 \\ 0 & 100 \\ 1 & 1 \end{pmatrix}$$

Computing  $\mathcal{Y} = \mathcal{T} \times_1 \mathbf{M}$  gives

$$\mathbf{Y}_1 = \begin{pmatrix} 10 & 30 \\ 200 & 400 \\ 3 & 7 \end{pmatrix} \quad \mathbf{Y}_2 = \begin{pmatrix} 50 & 60 \\ 600 & 800 \\ 11 & 15 \end{pmatrix}$$

# The Tucker3 Tensor Decomposition



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

# Tucker3 Decomposition

- The **Tucker3 tensor decomposition** decomposes the tensor into three **factor matrices  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$** , and a **core tensor  $\mathcal{G}$** 
  - $\mathbf{A}$  has  $P$ ,  $\mathbf{B}$  has  $Q$ , and  $\mathbf{C}$  has  $R$  columns and  $\mathcal{G}$  is  $P$ -by- $Q$ -by- $R$
- Many degrees of freedom: often  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  are required to be orthogonal
- If  $P=Q=R$  and core tensor  $\mathcal{G}$  is **hyper-diagonal**, then Tucker3 decomposition reduces to CP decomposition

# Solving Tucker3

- ALS-style methods are typically used

- The matricized forms are

$$\mathbf{X}_{(1)} = \mathbf{A}\mathbf{G}_{(1)}(\mathbf{C} \otimes \mathbf{B})^T$$

$$\mathbf{X}_{(2)} = \mathbf{B}\mathbf{G}_{(2)}(\mathbf{C} \otimes \mathbf{A})^T$$

$$\mathbf{X}_{(3)} = \mathbf{C}\mathbf{G}_{(3)}(\mathbf{B} \otimes \mathbf{A})^T$$

- If factor matrices are orthogonal, we can get  $\mathcal{G}$  as  $\mathcal{G} = \mathcal{X} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$

# **HOSVD, Tucker2, RESCAL, and DEDICOM**

- There are many tensor decompositions that are based on or similar to Tucker3
  - Or merge Tucker3 and CP
- Here are few, but the list is by no means exhaustive

# Higher-Order SVD (HOSVD)

- One method to compute the Tucker3 decomposition
  - Set **A** as the leading  $P$  left singular vectors of  $\mathbf{X}_{(1)}$
  - Set **B** as the leading  $Q$  left singular vectors of  $\mathbf{X}_{(2)}$
  - Set **C** as the leading  $R$  left singular vectors of  $\mathbf{X}_{(3)}$
  - Set tensor  $\mathcal{G}$  as  $\mathcal{X} \times_1 \mathbf{A}^T \times_2 \mathbf{B}^T \times_3 \mathbf{C}^T$

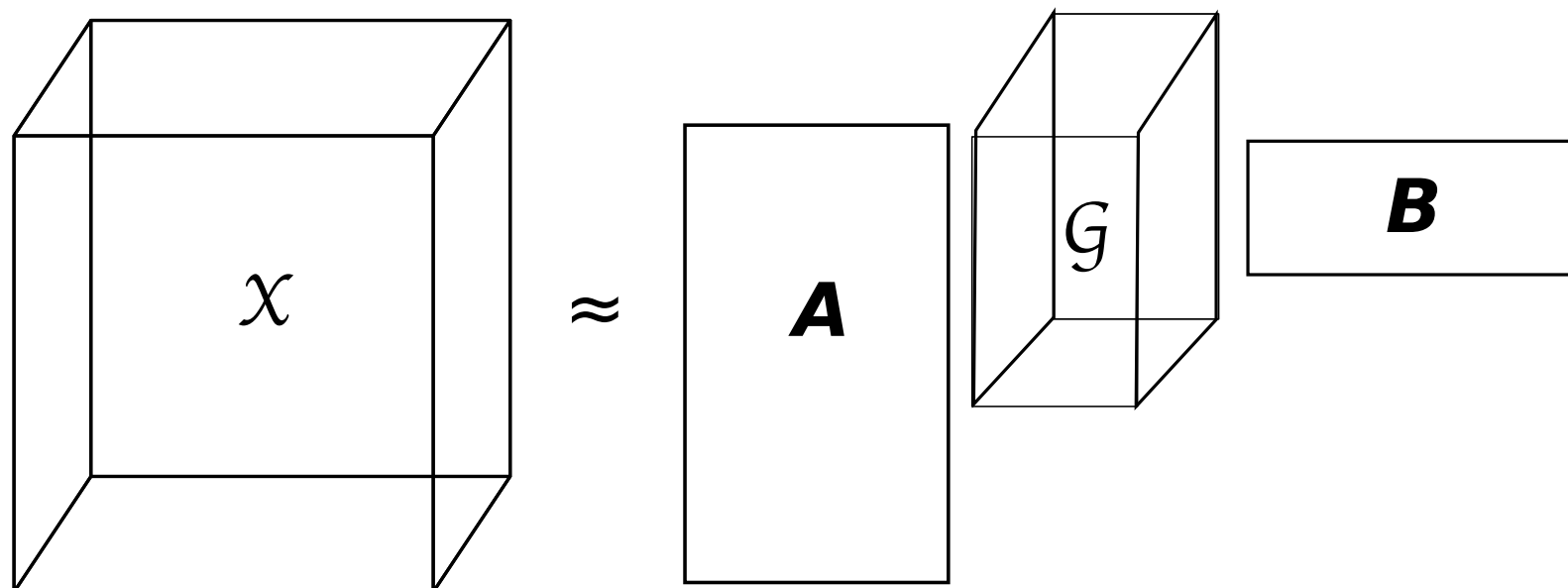


# Why HOSVD?

- Can be used as is for data analysis
  - E.g. TensorFaces
- Can be used to initialize other Tucker3 algorithms
  - Instead of random **A**, **B**, and **C**

# Tucker2 Decomposition

- The Tucker2 decomposition decomposes a 3-way tensor into a core tensor and two factor matrices
- Or, third factor matrix is forced to be an identity matrix
- Core keeps that mode's dimensionality



# Tucker2 Sliced and Matricized

- The slice-wise Tucker2:  $\mathbf{X}_k = \mathbf{A}\mathbf{G}_k\mathbf{B}^T$  for each  $k$
- Matricized forms replace  $\mathbf{C}$  with identity matrix  $\mathbf{I}$ :  $\mathbf{X}_{(1)} = \mathbf{A}\mathbf{G}_{(1)}(\mathbf{I} \otimes \mathbf{B})^T$  etc.
- To compute Tucker2:
  - Solve  $\mathbf{A}$  and  $\mathbf{B}$  using the matricized forms
  - Update each frontal slice of  $\mathbf{G}$  separately

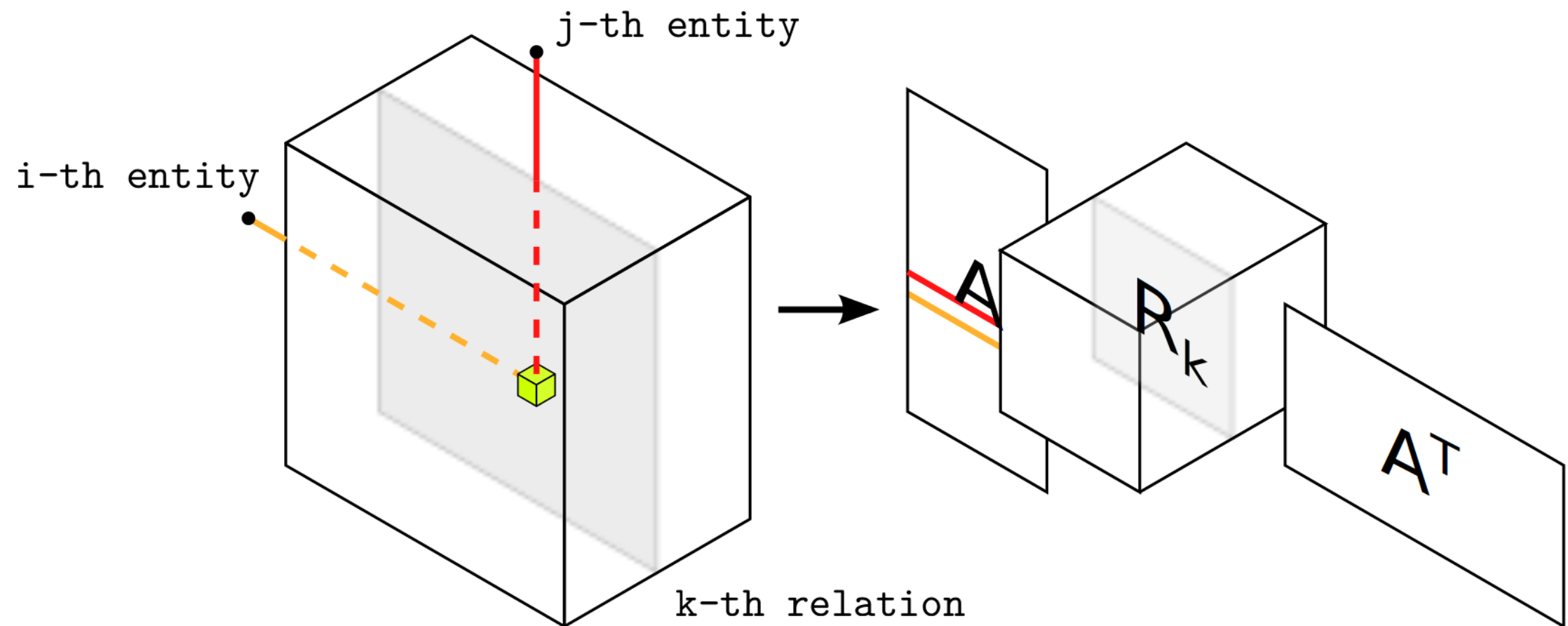
# Why Tucker2?

- Use Tucker2 if you don't want to factorize one mode
  - Too small dimension (e.g. 500-by-300-by-3)
  - This mode requires separate handling
    - E.g. if third mode is time, first Tucker2 and then time-series analysis on third mode
- Tucker2 is slightly simpler than Tucker3

# The RESCAL Decomposition

- The **RESCAL decomposition** merges Tucker2 and INDSCAL
- Tensor  $\mathcal{X}$  is factored into one factor matrix **A** and one core tensor **R**
  - $\mathbf{X}_k = \mathbf{A} \mathbf{R}_k \mathbf{A}^T$
- Tensor  $\mathcal{X}$  might not be symmetric on first two modes

# RESCAL in Picture



# Computing RESCAL (1)

- Mode-1 matricization of RESCAL is

$$\mathbf{X}_{(1)} = \mathbf{A} \mathbf{R}_{(1)} (\mathbf{I} \otimes \mathbf{A})^T$$

- This is hard as  $\mathbf{A}$  is both left and right
- Simplify: place pairs  $(\mathbf{X}_k \mathbf{X}_k^T)$  side-by-side and consider the right  $\mathbf{A}$  fixed
  - The  $\mathbf{X}_k^T$  guide  $\mathbf{A}$  to fit well also in RHS

# Computing RESCAL (2)

- To minimize the error, we minimize

$$||\mathbf{Y} - \mathbf{A}\mathbf{H}(\mathbf{I}_{2K} \otimes \mathbf{A}^T)||_F$$

- $\mathbf{Y} = [\mathbf{X}_1 \mathbf{X}_1^T \mathbf{X}_2 \mathbf{X}_2^T \dots \mathbf{X}_K \mathbf{X}_K^T]$

- $\mathbf{H} = [\mathbf{R}_1 \mathbf{R}_1^T \mathbf{R}_2 \mathbf{R}_2^T \dots \mathbf{R}_K \mathbf{R}_K^T]$

- For fixed  $\mathbf{A}^T$  and  $\mathcal{R}$ , the update rule for  $\mathbf{A}$  is

$$\mathbf{A} = \left( \sum_{k=1}^K (\mathbf{X}_k \mathbf{R}_k^T + \mathbf{X}_k^T \mathbf{A} \mathbf{R}_k) \right) \left( \sum_{k=1}^K (\mathbf{B}_k + \mathbf{C}_k) \right)^{-1}$$

- Here,  $\mathbf{B}_k = \mathbf{R}_k \mathbf{A}^T \mathbf{A} \mathbf{R}_k^T$  and  $\mathbf{C}_k = \mathbf{R}_k^T \mathbf{A}^T \mathbf{A} \mathbf{R}_k$



# Computing RESCAL (3)

- Each slice  $\mathbf{R}_k$  can be updated separately
  - Minimize  $||\text{vec}(\mathbf{X}_k) - (\mathbf{A} \otimes \mathbf{A})\text{vec}(\mathbf{R}_k)||$ 
    - Linear regression, set  $\text{vec}(\mathbf{R}_k) = (\mathbf{A} \otimes \mathbf{A})^+ \text{vec}(\mathbf{X}_k)$
  - To avoid computing the pseudo-inverse of big  $\mathbf{A} \otimes \mathbf{A}$ , compute the skinny QR decomposition of  $\mathbf{A}$ 
    - $\mathbf{A} = \mathbf{Q}\mathbf{U}$ ,  $\mathbf{Q}$  column-orthogonal,  $\mathbf{U}$  upper-triangular
    - Now:  $||\mathbf{X}_k - \mathbf{A}\mathbf{R}_k\mathbf{A}^T|| = ||\mathbf{X}_k - \mathbf{Q}\mathbf{U}\mathbf{R}_k\mathbf{U}^T\mathbf{Q}^T||$   
 $= ||\mathbf{Q}^T\mathbf{X}_k\mathbf{Q} - \mathbf{U}\mathbf{R}_k\mathbf{U}^T||$  and update rule as  $(\mathbf{U} \otimes \mathbf{U})$   
which is only  $R^2$ -by- $R^2$

# Why RESCAL

- No factorization of the third mode
  - Same as in Tucker2
- Only one factor matrix
  - We assume some kind of symmetry (INDSCAL)
    - E.g. subjects and objects
  - Provides "information flow" between the modes
- Each frontal slice has a separate "mixing matrix" for the interactions between factors

# The DEDICOM Decomposition: Matrix Version

- The **DEDICOM decomposition** is a matrix decomposition for an asymmetric relation between entities
  - What is the value of export from country  $i$  to country  $j$ ?
  - How many emails person  $i$  sent to person  $j$ ?
- **$X = ARA^T$** 
  - **$A$**  factors the entities
  - **$R$**  explains the asymmetric relation

# The DEDICOM Decomposition: Tensor Version

- The **three-way DEDICOM** adds weights for each factor's participation in each position in the third mode
  - E.g. how much country factor  $r$  acts as a seller or buyer at time  $k$ ?
- $\mathbf{X}_k = \mathbf{A}\mathbf{D}_k\mathbf{R}\mathbf{D}_k\mathbf{A}^T$ 
  - $\mathbf{A}$  and  $\mathbf{R}$  as before,  $\mathcal{D}$  is  $R$ -by- $R$ -by- $K$  tensor such that each frontal slice  $\mathbf{D}_k$  is diagonal
  - $(\mathbf{D}_k)_{rr}$  is the weight for factor  $r$  at time  $k$

# DEDICOM in Picture

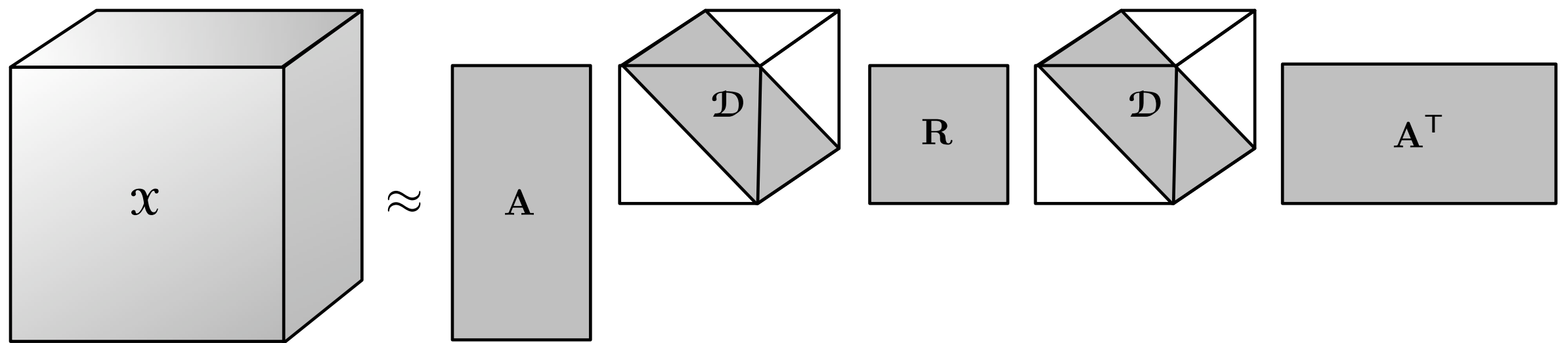


Fig. 5.2: Three-way DEDICOM model.

# Computing DEDICOM: ASALSAN (1)

- We want to minimize  $\sum_k ||\mathbf{X}_k - \mathbf{A}\mathbf{D}_k\mathbf{R}\mathbf{D}_k\mathbf{A}^T||$
- ASALSAN (Alternating Simultaneous Approximation, Least Squares, and Newton) is one way
  - Stack pairs  $(\mathbf{X}_k \mathbf{X}_k^T)$ :  $\mathbf{Y} = [\mathbf{X}_1 \mathbf{X}_1^T \dots \mathbf{X}_K \mathbf{X}_K^T]$
  - We get  $||\mathbf{Y} - \mathbf{A}\mathbf{H}(\mathbf{I}_{2K} \otimes \mathbf{A}^T)||$  with  
 $\mathbf{H} = [\mathbf{D}_1\mathbf{R}\mathbf{D}_1 \mathbf{D}_1\mathbf{R}^T\mathbf{D}_1 \dots \mathbf{D}_K\mathbf{R}\mathbf{D}_K \mathbf{D}_K\mathbf{R}^T\mathbf{D}_K]$

# Computing DEDICOM: ASALSAN (2)

- To update  $\mathbf{A}$ , fix right  $\mathbf{A}$  and update the left  

$$\mathbf{A} = \left( \sum_{k=1}^K (\mathbf{X}_k \mathbf{A} \mathbf{D}_k \mathbf{R}^T \mathbf{D}_k + \mathbf{X}_k^T \mathbf{A} \mathbf{D}_k \mathbf{R} \mathbf{D}_k) \right) \left( \sum_{k=1}^K (\mathbf{B}_k + \mathbf{C}_k) \right)^{-1}$$
  - $\mathbf{B}_k = \mathbf{D}_k \mathbf{R} \mathbf{D}_k (\mathbf{A}^T \mathbf{A}) \mathbf{D}_k \mathbf{R}^T \mathbf{D}_k$  and  
 $\mathbf{C}_k = \mathbf{D}_k \mathbf{R}^T \mathbf{D}_k (\mathbf{A}^T \mathbf{A}) \mathbf{D}_k \mathbf{R} \mathbf{D}_k$
- To update  $\mathbf{R}$ , we use vectors:  

$$\min_{\mathbf{R}} \left\| \begin{pmatrix} \text{vec}(\mathbf{X}_1) \\ \vdots \\ \text{vec}(\mathbf{X}_K) \end{pmatrix} - \begin{pmatrix} \mathbf{A} \mathbf{D}_1 \otimes \mathbf{A} \mathbf{D}_1 \\ \vdots \\ \mathbf{A} \mathbf{D}_K \otimes \mathbf{A} \mathbf{D}_K \end{pmatrix} \text{vec}(\mathbf{R}) \right\|$$
- To update  $\mathcal{D}$ , use Newton's method for each slice  $\mathbf{D}_k$

# DEDICOM vs. RESCAL vs. INDSCAL vs. Tucker2

- RESCAL is a relaxed version of DEDICOM
  - Mixing matrix  **$R$**  is different for each slice
  - Easier to compute as there's no tensor  $\mathcal{D}$ 
    - Algorithm similar to ASALSAN, but simpler
- RESCAL is to Tucker2 what INDSCAL is to CP
  - Share's INDSCAL's equal factor matrix
  - Uses Tucker2's core



# The Non-Negative Variants

- Sometimes having non-negative factors is beneficial for data analysis
  - Improved interpretability
    - E.g. physical measurements
  - Sparsity
- All of the discussed methods can be cast into non-negative variants

# Non-Negative CP

- The simplest way to compute non-negative CP is to use non-negative least-squares solver with the matricized equations
  - $\min_{\mathbf{A} \in \mathbb{R}_+^{N \times R}} \left\| \mathbf{X}_{(1)} - \mathbf{A}(\mathbf{C} \odot \mathbf{B})^T \right\|$
- Also multiplicative updates are possible
  - $a_{ir} = a_{ir} \frac{(\mathbf{X}_{(1)} \mathbf{Z})_{ir}}{(\mathbf{A} \mathbf{Z}^T \mathbf{Z})_{ir}}$  with  $\mathbf{Z} = (\mathbf{C} \odot \mathbf{B})$
- Also other methods exist

# Non-Negative Others

- Also non-negative Tucker[2|3] can be solved using multiplicative update rules
- Non-negative ASALSAN yields non-negative DEDICOM
  - Similar algorithm will work for RESCAL

# Summary

- Many, many different tensor decomposition
  - User's responsibility to choose the correct one
    - How do the results look like?
    - What's the time complexity?
- Most algorithms are geared towards dense data
  - But many data analysis data are sparse

# Suggested Reading

- In addition to those from last time:
- Acar, E., & Yener, B. (2009). Unsupervised Multiway Data Analysis: A Literature Survey. *IEEE Transactions on Knowledge and Data Engineering*, 21(1), 6–20. doi:10.1109/TKDE.2008.112
- Shorter and more focused on applications than Kolda & Bader (2009)