

Tensor Factorizations & Applications

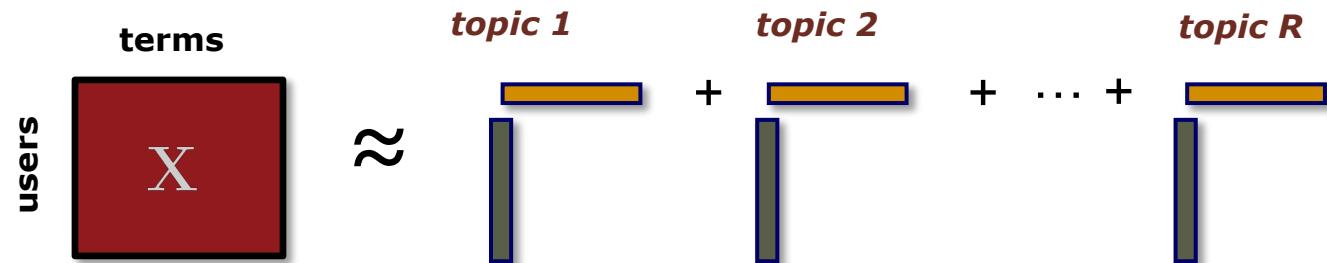
Evrim Acar

**Simula Metropolitan Center for Digital Engineering
Oslo, Norway**

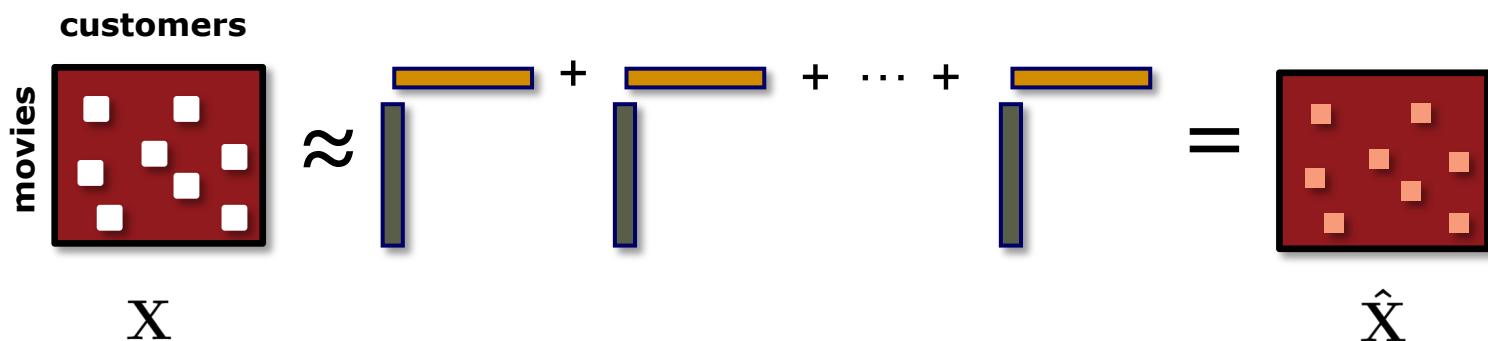
simula

Matrix Factorizations in Data Mining

Capturing underlying hidden factors/patterns



Incomplete Matrix Factorization/Matrix Completion



simula

Data is often squeezed to be two-way!

2-way

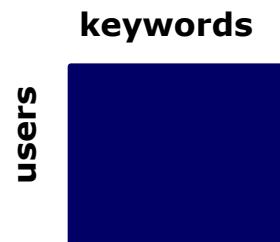
Social networks : **<users, keywords>**
<employee, employee>
<users, users>

Text mining: **< documents, terms>**

Recommender
Systems: **<customers, items>**

Computer vision: **<people, pixels>**

Neuroscience: **<electrodes, time>**



Data is often squeezed to be two-way!

2-way

Social networks : <users, keywords>
<employee, employee>
<users, users>

Text mining: < documents, terms>

Recommender
Systems: <customers, items>

Computer vision: <people, pixels>

Neuroscience: <electrodes, time>

Multi-way

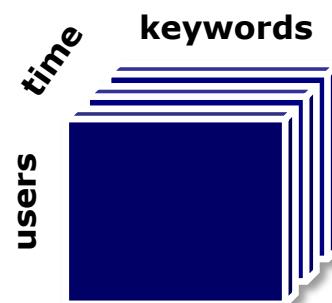
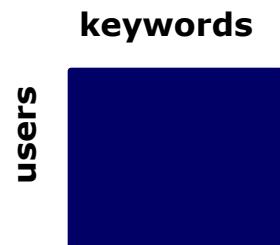
< users, keywords, time>
< employee, employee, time>
< users, users, communication type>

< documents, terms, terms>

<customers, items, tags>

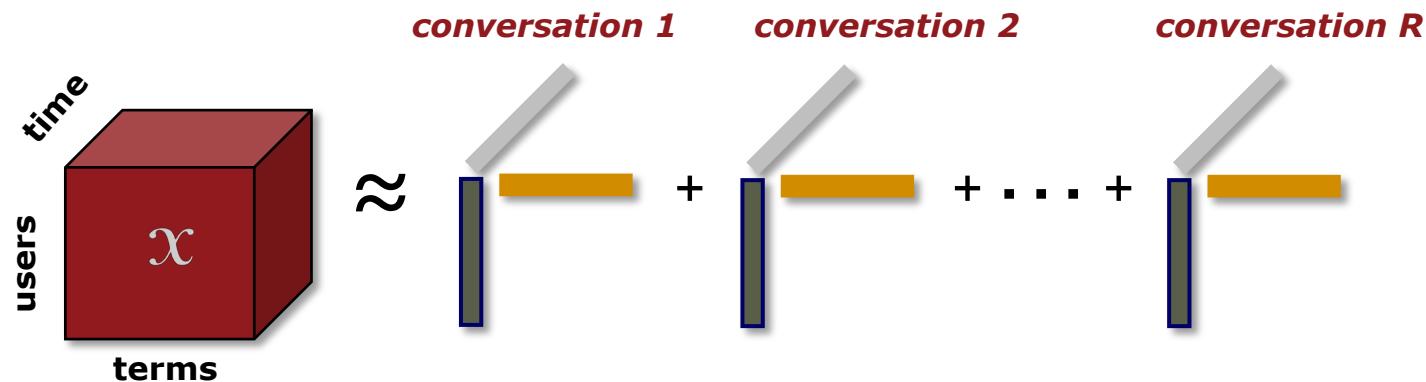
<people, pixels, viewpoints>

<electrodes, time, frequency>
<electrodes, time, subjects>

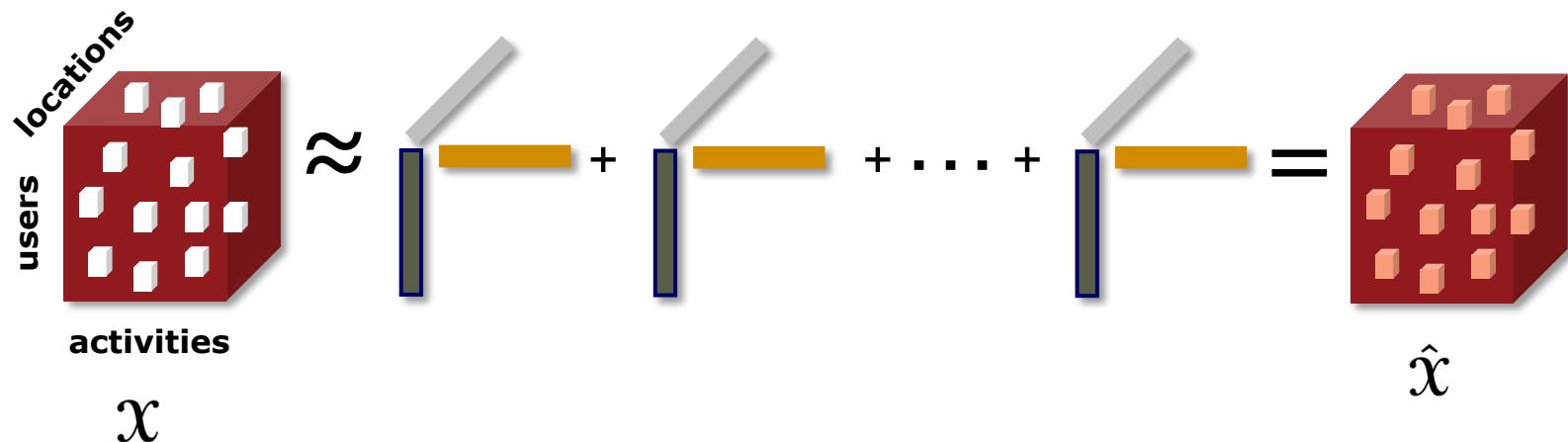


Data is often multi-way!

Capturing underlying hidden factors/patterns



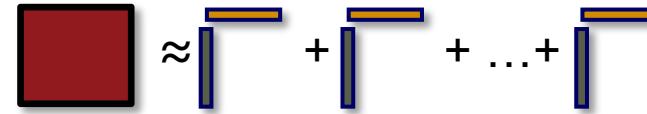
Incomplete Tensor Factorization/
Tensor Completion



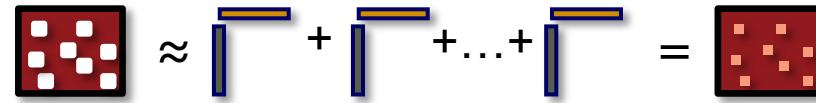
simula

Matrix Factorizations in Data Mining

Matrix Factorizations

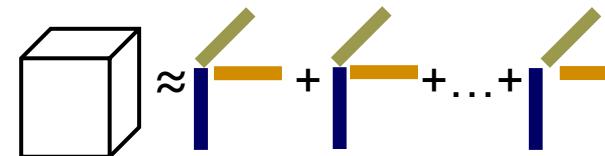


Matrix Completion

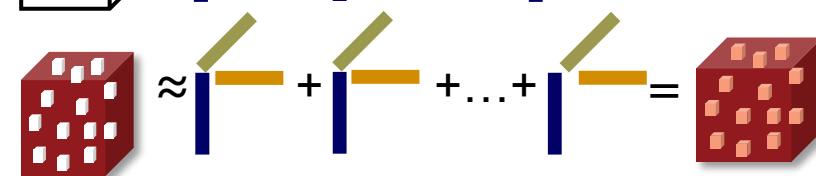


Data sets are often multi-way: Tensor Factorizations

Tensor Factorizations

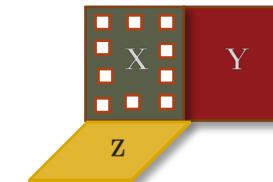


Tensor Completion

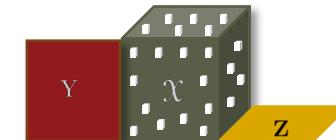
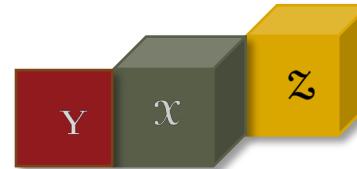


Data sets often come from multiple sources: Data Fusion

Coupled Matrix Factorizations

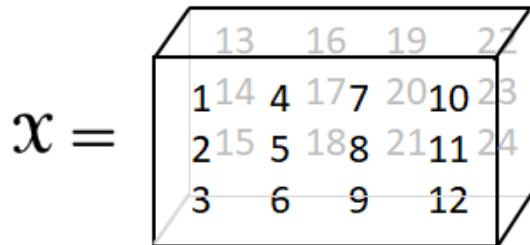
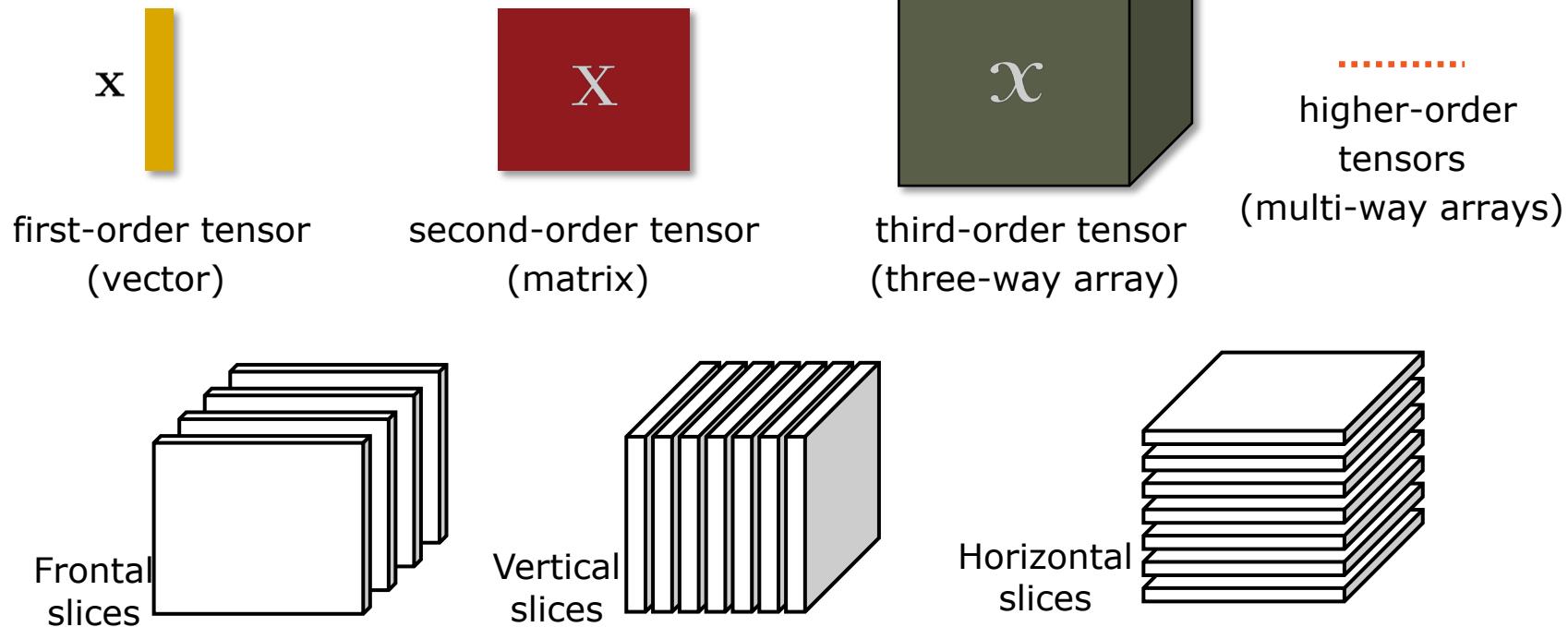


Coupled Tensor Factorizations



Tensors: Terminology

Order: Number of Modes or Ways



Order: 3

Size: $3 \times 4 \times 2$

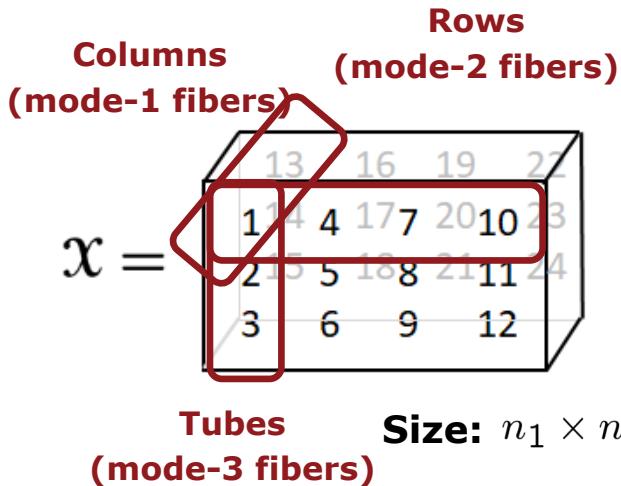
$$X(:,:,1) = \begin{bmatrix} 1 & 4 & 7 & 10 \\ 2 & 5 & 8 & 11 \\ 3 & 6 & 9 & 12 \end{bmatrix}$$

Frontal slices

$$X(:,:,2) = \begin{bmatrix} 13 & 16 & 19 & 22 \\ 14 & 17 & 20 & 23 \\ 15 & 18 & 21 & 24 \end{bmatrix}$$

Tensors: Basic Operations

Matricization (a.k.a. Unfolding or Flattening)



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

$n_1 \times n_2 n_3$

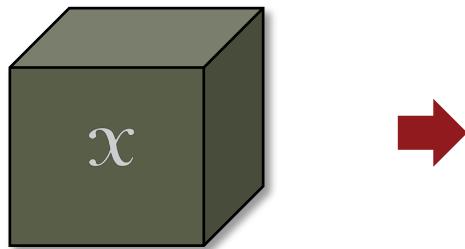
$$\mathbf{X}_{(2)} = \begin{bmatrix} 1 & 2 & 3 & 13 & \dots & 15 \\ 4 & 5 & 6 & 16 & \dots & 18 \\ 7 & 8 & 9 & 19 & \dots & 21 \\ 10 & 11 & 12 & 22 & \dots & 24 \end{bmatrix}$$

$n_2 \times n_1 n_3$

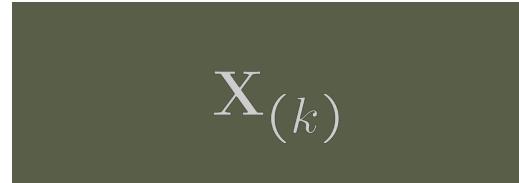
$$\mathbf{X}_{(3)} = \begin{bmatrix} 1 & 2 & 3 & 4 & \dots & 12 \\ 13 & 14 & 15 & 16 & \dots & 24 \end{bmatrix}$$

$n_3 \times n_1 n_2$

Unfolding a d -way array (mode- k unfolding)



Size: $n_1 \times n_2 \times \dots \times n_d$



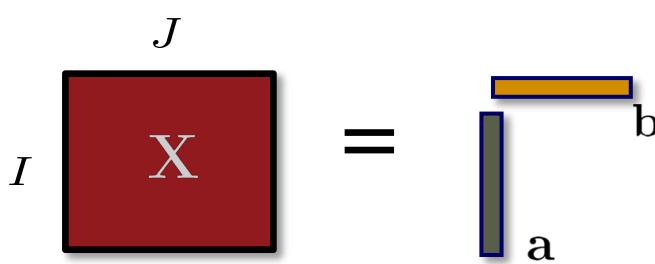
$$n_k \times \prod_{i,i \neq k}^d n_i$$

simula

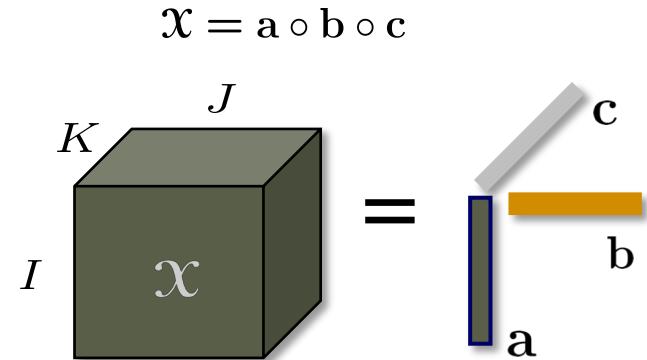
Tensors: Basic Operations (cont.)

Vector outer product $\mathbf{a} \in \mathbb{R}^I, \mathbf{b} \in \mathbb{R}^J, \mathbf{c} \in \mathbb{R}^K$

$$\mathbf{X} = \mathbf{ab}^\top = \mathbf{a} \circ \mathbf{b}$$



$$\mathbf{X} \in \mathbb{R}^{I \times J}, x_{ij} = a_i b_j$$



$$\mathbf{x} \in \mathbb{R}^{I \times J \times K}, x_{ijk} = a_i b_j c_k$$

Kronecker product

$\mathbf{a} \in \mathbb{R}^I, \mathbf{b} \in \mathbb{R}^J, \mathbf{A} \in \mathbb{R}^{I \times M}, \mathbf{B} \in \mathbb{R}^{J \times N}$

$$\mathbf{a} \otimes \mathbf{b} = \begin{bmatrix} a_1 \mathbf{b} \\ a_2 \mathbf{b} \\ \vdots \\ a_I \mathbf{b} \end{bmatrix} \in \mathbb{R}^{IJ}$$

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11} \mathbf{B} & \dots & a_{1M} \mathbf{B} \\ a_{21} \mathbf{B} & \dots & a_{2M} \mathbf{B} \\ \vdots & & \vdots \\ a_{I1} \mathbf{B} & \dots & a_{IM} \mathbf{B} \end{bmatrix} \in \mathbb{R}^{IJ \times MN}$$

Khatri-Rao product (Columnwise Kronecker Product)

$\mathbf{A} \in \mathbb{R}^{I \times R}, \mathbf{B} \in \mathbb{R}^{J \times R}$

$$\mathbf{A} \odot \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \mathbf{a}_2 \otimes \mathbf{b}_2 \quad \dots \quad \mathbf{a}_R \otimes \mathbf{b}_R] \in \mathbb{R}^{IJ \times R}$$

Tensors: Basic Operations (cont.)

Tensor Inner Product $\mathcal{X} \in \mathbb{R}^{I \times J \times K}, \mathcal{Y} \in \mathbb{R}^{I \times J \times K}$

$$\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{k=1}^K \sum_{j=1}^J \sum_{i=1}^I x_{ijk} y_{ijk} \quad \langle \begin{matrix} \mathcal{X} \\ \mathcal{Y} \end{matrix} \rangle = c$$

Frobenius Norm

$$\|\mathcal{X}\|^2 = \langle \mathcal{X}, \mathcal{X} \rangle = \sum_{k=1}^K \sum_{j=1}^J \sum_{i=1}^I x_{ijk}^2$$

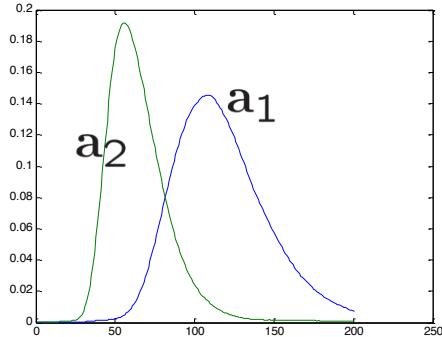
Hadamard Product $\mathcal{X}, \mathcal{Y}, \mathcal{Z} \in \mathbb{R}^{I \times J \times K}$

$$\mathcal{Z} = \mathcal{X} * \mathcal{Y} \quad z_{ijk} = x_{ijk} y_{ijk}$$

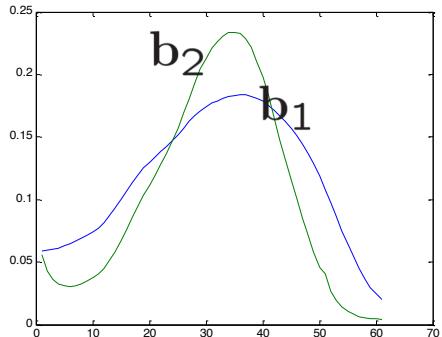
What if we have multiple matrices with the same underlying factors but in different proportions...

True factors

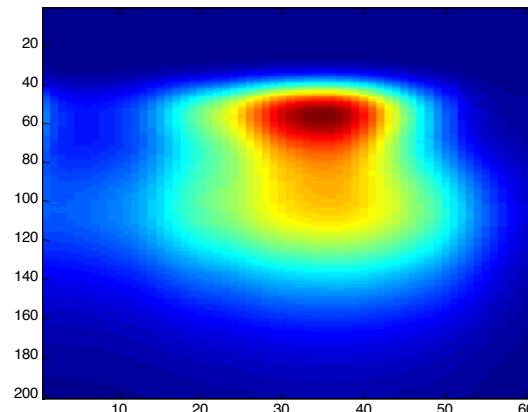
$$A \in \mathbb{R}^{201 \times 2} = [a_1 \ a_2]$$



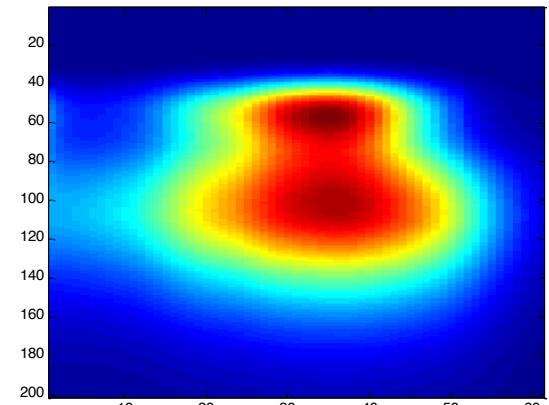
$$B \in \mathbb{R}^{61 \times 2} = [b_1 \ b_2]$$



$$X_1 = A \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} B^T$$



$$X_2 = A \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} B^T$$



We can recover the true factors uniquely up to trivial indeterminacies, i.e., scaling and permutation.

Canonical Polyadic (CP) CANDECOMP/PARAFAC (CP)

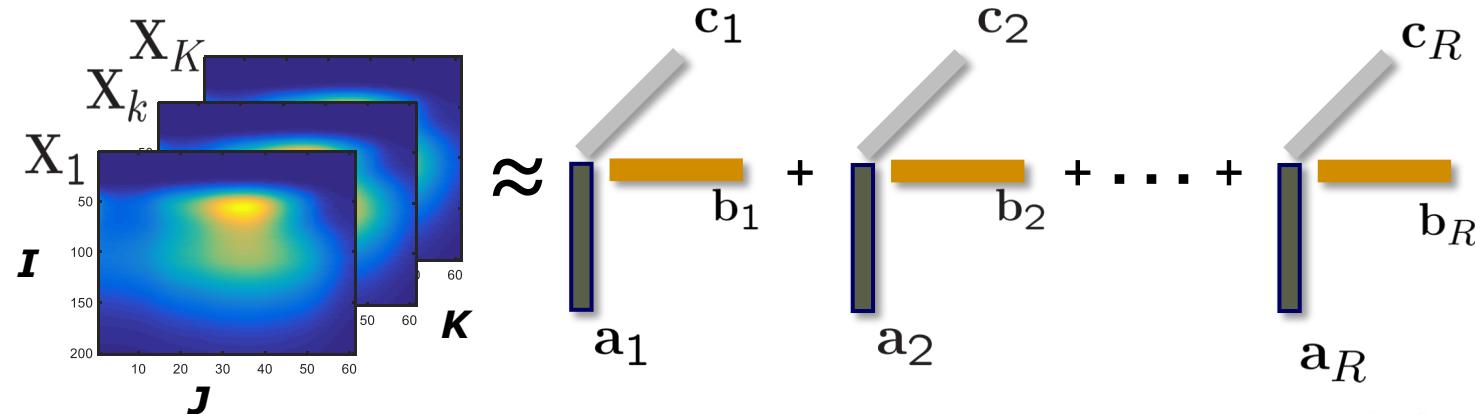
Hitchcock, 1927: Polyadic form of a tensor

Harshman, 1970: Parallel Factor Analysis (PARAFAC)

Carroll & Chang, 1970: Canonical Decomposition (CANDECOMP)

A popular tensor factorization model: CP

As an extension of matrix factorizations to higher-order tensors (multi-way arrays), tensor factorizations are used to extract the underlying factors in higher-order data sets. The CP model represents a tensor as a sum of rank-one tensors:



$$X_k \approx A \begin{bmatrix} c_{k1} & \dots & c_{kR} \end{bmatrix} B^T$$

$$\begin{aligned} \mathcal{X} &\approx \sum_{r=1}^R a_r \circ b_r \circ c_r \\ &\approx [\![A, B, C]\!] \end{aligned}$$

$$A \in \mathbb{R}^{I \times R} = [a_1 \ \dots \ a_R]$$

$$B \in \mathbb{R}^{J \times R} = [b_1 \ \dots \ b_R]$$

$$C \in \mathbb{R}^{K \times R} = [c_1 \ \dots \ c_R]$$

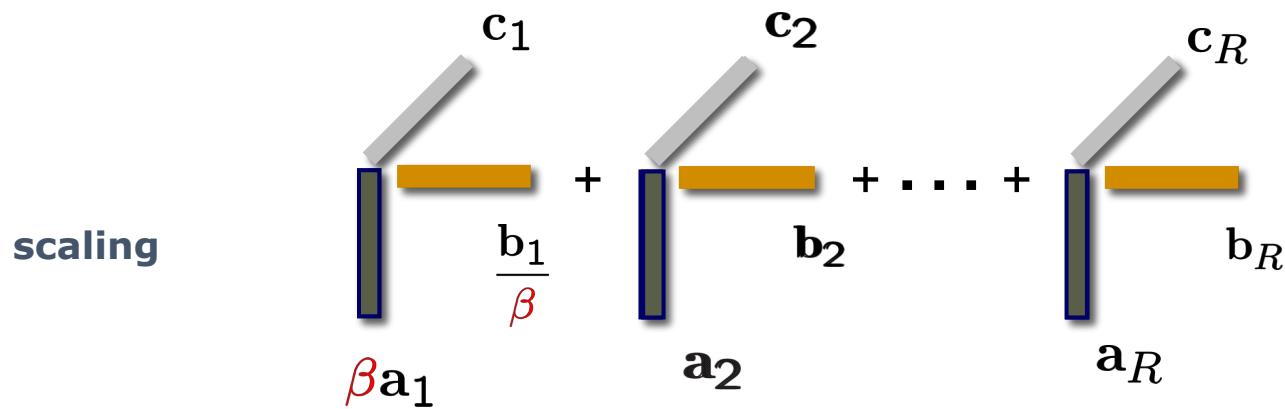
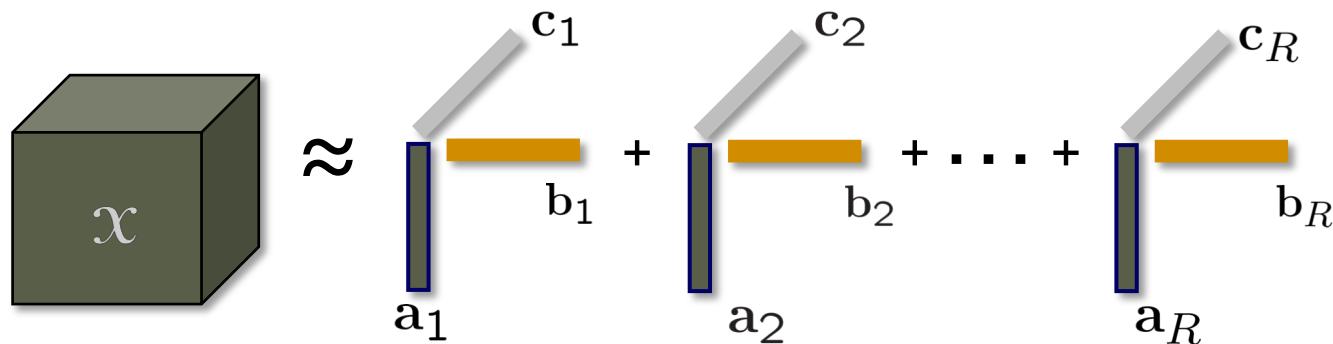
CP is unique up to scaling and permutation under the condition that

$$\text{krank}(A) + \text{krank}(B) + \text{krank}(C) \geq 2R + 2$$

[Kruskal, 1977;
Sidiropoulos and Bro, 2000]

where krnk (A) = max. value of k such that any k columns of A is linearly independent.

CP is unique up to trivial indeterminacies



permutation

+ . . . +

SVD vs. CP

SVD expresses a matrix as the sum of rank-one matrices.

$$\begin{matrix} X \\ \boxed{X} \end{matrix} = \sigma_1 u_1 v_1^\top + \sigma_2 u_2 v_2^\top + \dots + \sigma_R u_R v_R^\top$$
$$X = \sum_{r=1}^R \sigma_r u_r \circ v_r$$

CP expresses a higher-order tensor as the sum of rank-one tensors.

$$\begin{matrix} \mathcal{X} \\ \boxed{\mathcal{X}} \end{matrix} = \mathbf{a}_1 \mathbf{b}_1 \mathbf{c}_1^\top + \mathbf{a}_2 \mathbf{b}_2 \mathbf{c}_2^\top + \dots + \mathbf{a}_R \mathbf{b}_R \mathbf{c}_R^\top$$
$$\mathcal{X} = \sum_{r=1}^R \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

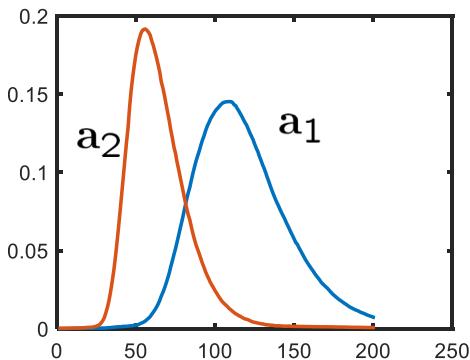
CP is an extension of SVD to higher-order tensors but it is different in many ways.

Does it really work?

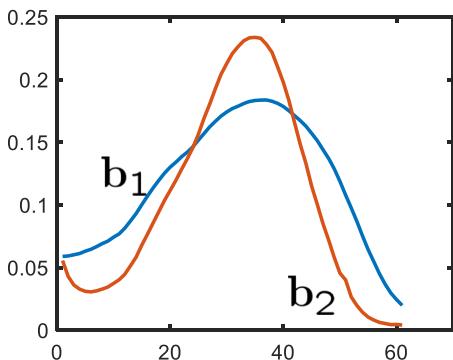
EXAMPLE:

True factors

$$A \in \mathbb{R}^{201 \times 2} = [a_1 \ a_2]$$

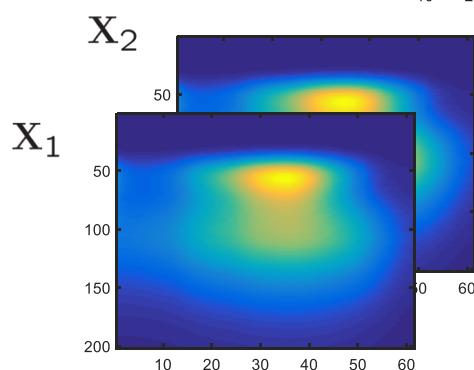
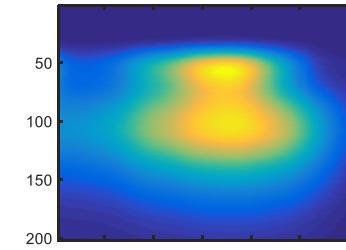
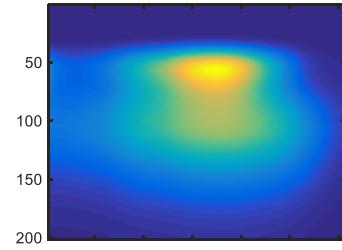


$$B \in \mathbb{R}^{61 \times 2} = [b_1 \ b_2]$$



$$X_1 = A \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} B^T$$

$$X_2 = A \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} B^T$$



$\mathcal{X} \approx [[A, B, C]]$

```
load toydata.mat
XX(:,:,1) = A*[1 0; 0 1]*B';
XX(:,:,2) = A*[3 0; 0 2]*B';
%set optimization parameters
options            = ncg('defaults');
options.DisplayIters = 50;
options.RelFuncTol = 1e-8;
options.StopTol    = 1e-8;
% call cp_opt to fit the CP model
[Fac,FacInit,out] = cp_opt(tensor(XX), 2,'init','nvecs','alg_options',options);
subplot(2,1,1);plot(Fac.U{1});
subplot(2,1,2);plot(Fac.U{2});
```

Algorithms for fitting the CP model

$$\mathcal{X} = \underbrace{\begin{array}{c} c_1 \\ \vdots \\ c_R \end{array}}_{\mathbf{C}} + \dots + \underbrace{\begin{array}{c} b_1 \\ \vdots \\ b_R \end{array}}_{\mathbf{B}} + \underbrace{\begin{array}{c} a_1 \\ \vdots \\ a_R \end{array}}_{\mathbf{A}}$$

The goal is to find matrices \mathbf{A} , \mathbf{B} , \mathbf{C} that solve the following optimization problem:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]\|^2$$

**Traditional Approach:
Alternating Least Squares (ALS)**
[Harshman, 1970; Carroll & Chang, 1970]

$$\begin{aligned}\mathbf{A} &= \mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})((\mathbf{C} \odot \mathbf{B})^T(\mathbf{C} \odot \mathbf{B}))^\dagger \\ &= \mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B})(\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B})^\dagger\end{aligned}$$

Matricized Tensor Times
Khatri-Rao Product (MTTKRP)

while “not converged” **do**

Solve for \mathbf{A} (with fixed \mathbf{B} , \mathbf{C})

$$\min_{\mathbf{A}} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]\|^2$$

Solve for \mathbf{B} (with fixed \mathbf{A} and \mathbf{C})

$$\min_{\mathbf{B}} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]\|^2$$

Solve for \mathbf{C} (with fixed \mathbf{A} and \mathbf{B})

$$\min_{\mathbf{C}} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]\|^2$$

end while

However, better convergence properties and accuracy have been achieved using all-at-once optimization approaches:

First-order methods: [Paatero, 1999], CP-OPT [Acar, Dunlavy, Kolda, 2011]

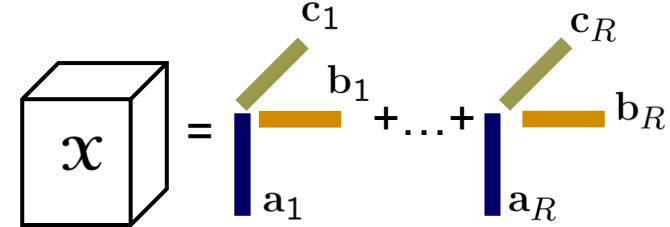
Second-order methods: [Paatero, 1997; Tomasi and Bro, 2006; Phan et al., 2013; Sorber et al., 2013]

CP-OPT: Gradient-based All-at-once Approach

[Acar, Dunlavy, and Kolda, 2011]

CP-OPT is a general gradient-based optimization approach for fitting a CP model.

$$\min_{A, B, C} \| \mathcal{X} - [A, B, C] \|^2$$



Step1: Define the objective function

$$f(A, B, C) = \frac{1}{2} \| \mathcal{X} - [A, B, C] \|^2$$

Step2: Compute the gradient *Remember the ALS steps!*

$$\frac{\partial f}{\partial A} = -X_{(1)}(C \odot B) + A(C^T C * B^T B)$$

$$\frac{\partial f}{\partial B} = -X_{(2)}(C \odot A) + B(C^T C * A^T A)$$

$$\frac{\partial f}{\partial C} = -X_{(3)}(B \odot A) + C(B^T B * A^T A)$$

→

Vectorize and concatenate the partials

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial a_1} \\ \vdots \\ \frac{\partial f}{\partial a_R} \\ \frac{\partial f}{\partial b_1} \\ \vdots \\ \frac{\partial f}{\partial b_R} \\ \frac{\partial f}{\partial c_1} \\ \vdots \\ \frac{\partial f}{\partial c_R} \end{bmatrix}$$

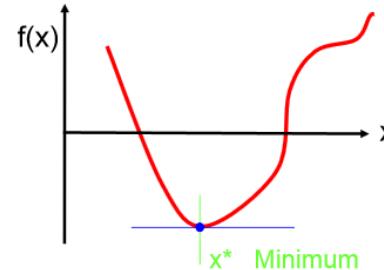
Step3: Pick a first-order optimization method

e.g., Nonlinear Conjugate Gradient (NCG) from the **Poblano Toolbox** [Dunlavy, Kolda and Acar, 2010]

Few words on optimization

Optimization problem

$$\min_x f(x)$$



Figures from Moritz Diehl's slides

Aim of most optimization algorithms

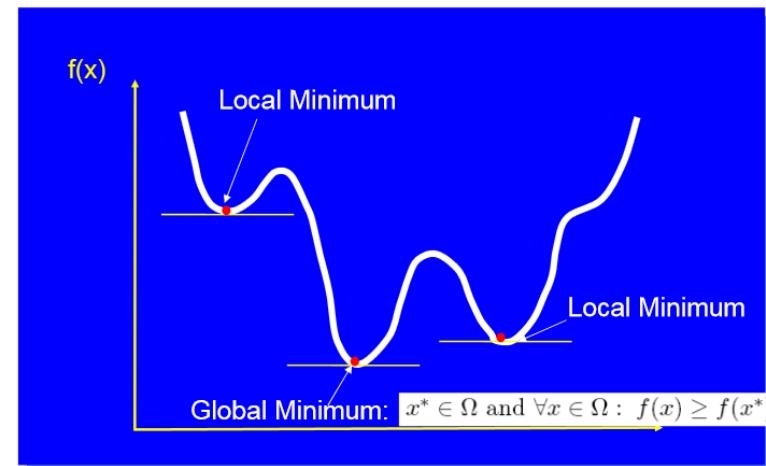
Find a local minimizer x^*

Gradient

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Hessian

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \ddots & \vdots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix}$$



First-order Necessary Conditions

If x^* is a local minimizer and f is continuously differentiable in an open neighborhood of x^* , then $\nabla f(x^*) = 0$.

Second-order Necessary Conditions

If x^* is a local minimizer of f and $\nabla^2 f$ exists and is continuous in an open neighborhood of x^* , then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.

Few words on optimization (cont.)

The general structure of the optimization algorithms in the Poblano Toolbox:

```
choose a starting value  $x_0$ 
while not converged do
    find the search direction  $p_k$ 
    determine the step length  $\alpha$ 
    new iterate  $x_k = x_{k-1} + \alpha p_k$ 
end while
```

Different choices of search direction and step length lead to different optimization algorithms.

Taylor's approximation:

Suppose that f is twice continuously differentiable, then,

$$f(x_k + p) \approx f_k + p^\top \nabla f_k + \frac{1}{2} p^\top \nabla^2 f_k p$$

Steepest descent:

$$p_k = -\nabla f_k$$

Newton's direction:

$$p_k = -(\nabla^2 f_k)^{-1} \nabla f_k$$

Quasi-Newton Methods (e.g., Limited-Memory BFGS):

$$p_k = -B_k^{-1} \nabla f_k$$

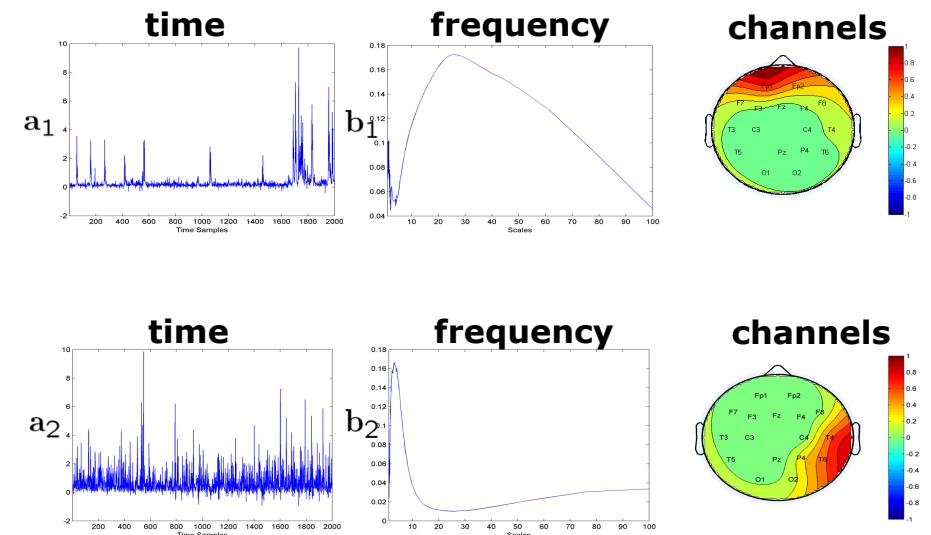
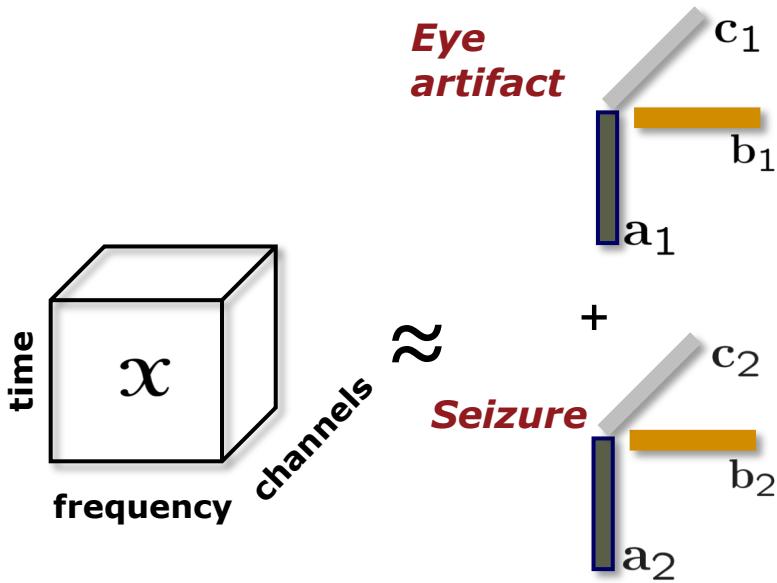
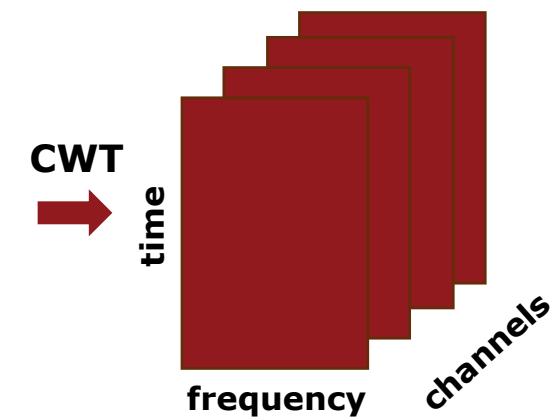
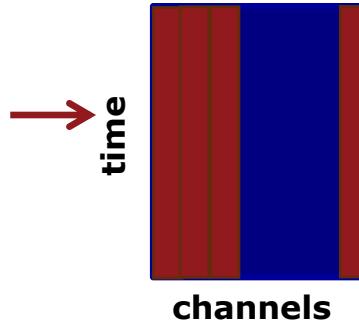
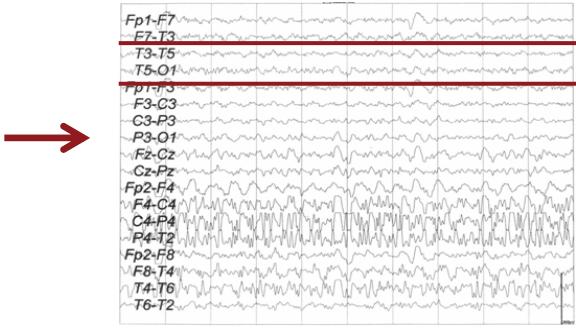
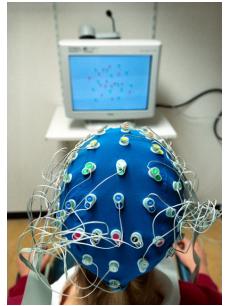
Nonlinear Conjugate Gradient Methods:

$$p_k = -\nabla f_k + \beta_k p_{k-1}$$

The ones using only the first-order info are implemented in the Poblano Toolbox.

Application₁ (Neuroscience): CP has proved useful in epileptic seizure localization

[Acar et al., 2007; De Vos et al., 2007]



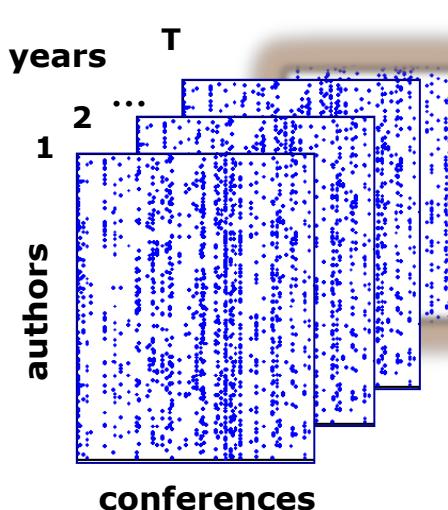
Application₂ (Recommender Systems): CP can capture temporal patterns useful for link prediction

[Dunlavy, Kolda, Acar, 2011]

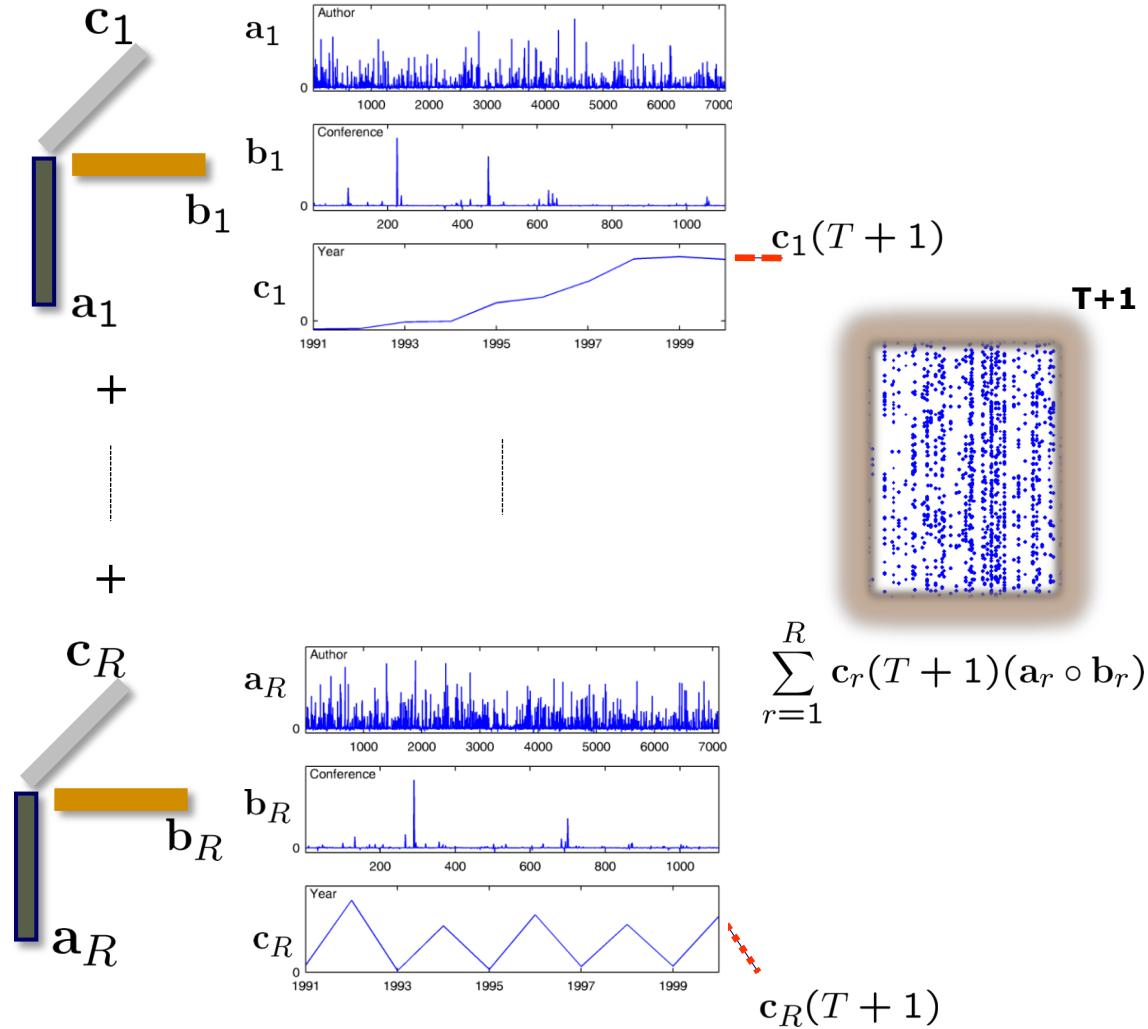
Temporal Link Prediction

Which customers will buy which products?

Which webpages will users visit?



2

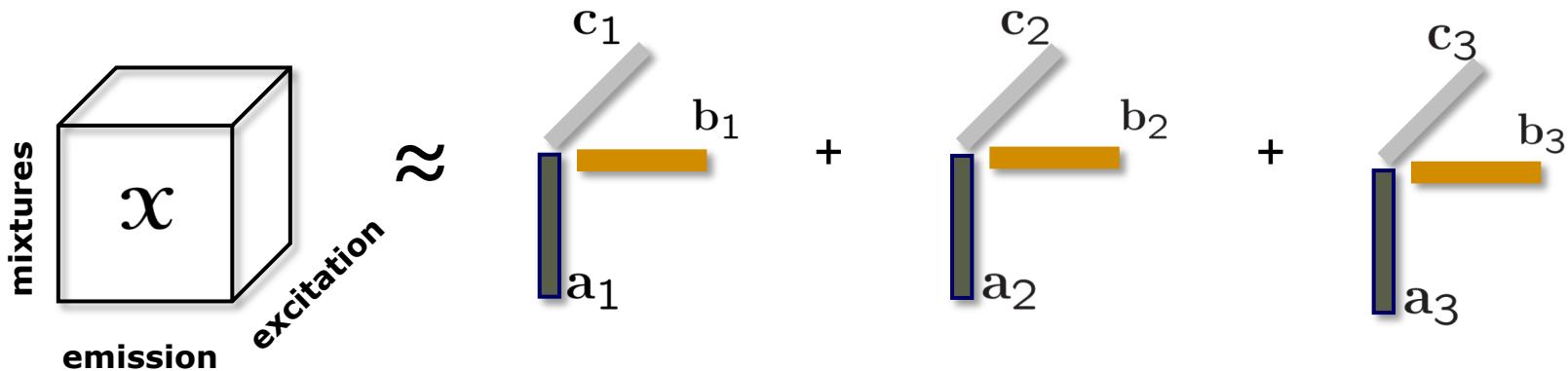


simula

Application₃ (Chemometrics): CP can separate the chemicals in mixtures

[Andersen and Bro, 2003]

A popular application of the CP model is the separation of individual chemicals from mixtures of chemicals measured using fluorescence spectroscopy.

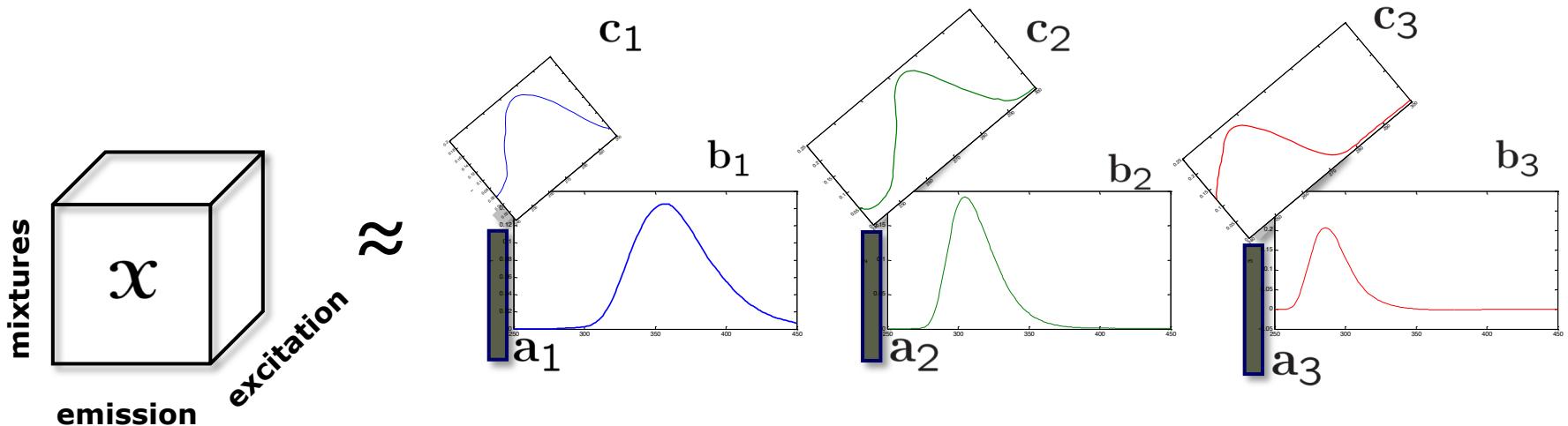


*Amino Acid Data
<http://www.models.life.ku.dk/>*

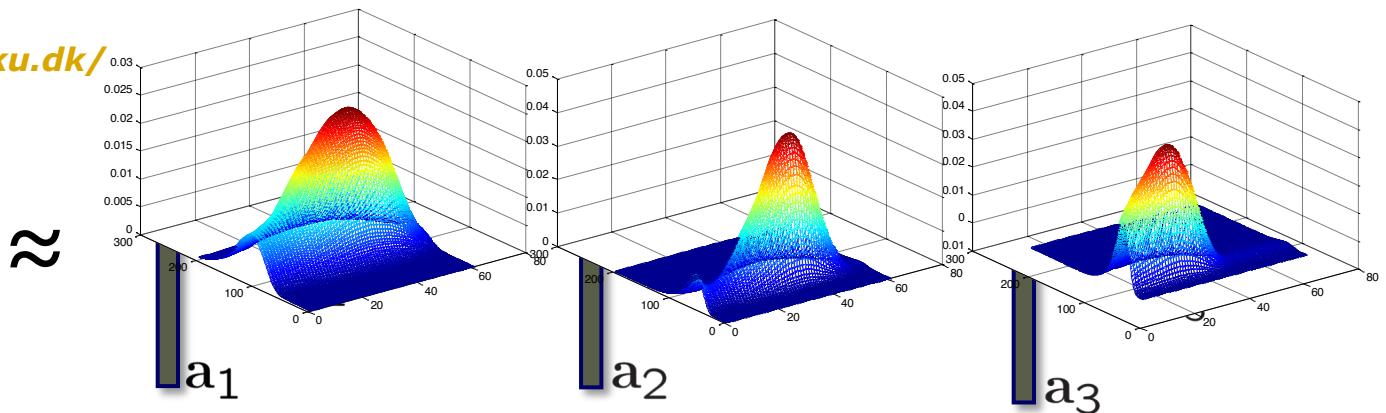
Application₃ (Chemometrics): CP can separate the chemicals in mixtures

[Andersen and Bro, 2003]

A popular application of the CP model is the separation of individual chemicals from mixtures of chemicals measured using fluorescence spectroscopy.



Amino Acid Data
<http://www.models.life.ku.dk/>



Modeling the aminoacids data

```
% Load data
load aminoacids.mat
X = tensor(X.data);
X = X/norm(X);
% Set optimization parameters and fit a CP model using CP-OPT
options            = ncg('defaults');
options.DisplayIters = 50;
options.RelFuncTol  = 1e-10;
options.StopTol     = 1e-10;
[Fac,FacInit,out]   = cp_opt(X, 3,'init','nvecs','alg','ncg','alg_options',options);
for i=1:3
    subplot(3,1,i)
    plot(Fac.U{i})
end
```

Model: CANDECOMP/PARAFAC (CP)

Algorithm: We can pick a gradient-based optimization method

Nonlinear Conjugate Gradient (ncg)

Limited-memory BFGS (lbfsgs)

....

Initialization of factor matrices:

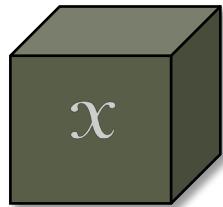
svd-based

entries randomly sampled from a standard normal distribution

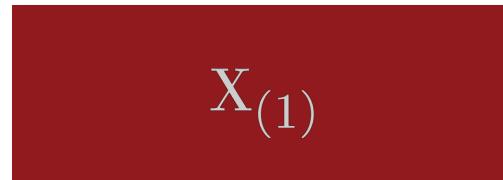
pre-defined cell array containing the factor matrices

SVD-based Initialization

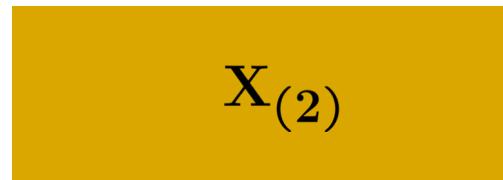
Left singular vectors (i.e., \mathbf{U}_1 , \mathbf{U}_2 , \mathbf{U}_3) of matricized tensors in each mode can be used to initialize the factor matrices.



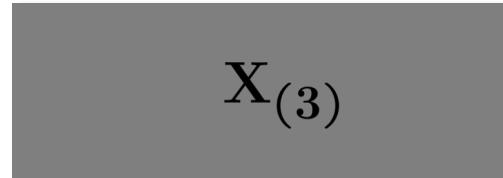
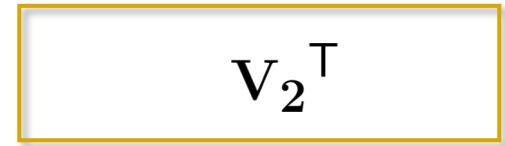
$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \| \mathcal{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \|^2$$



$$\approx \mathbf{U}_1 \begin{matrix} \diagdown \\ \Sigma_1 \end{matrix} \mathbf{V}_1^\top$$



$$\approx \mathbf{U}_2 \begin{matrix} \diagdown \\ \Sigma_2 \end{matrix} \mathbf{V}_2^\top$$

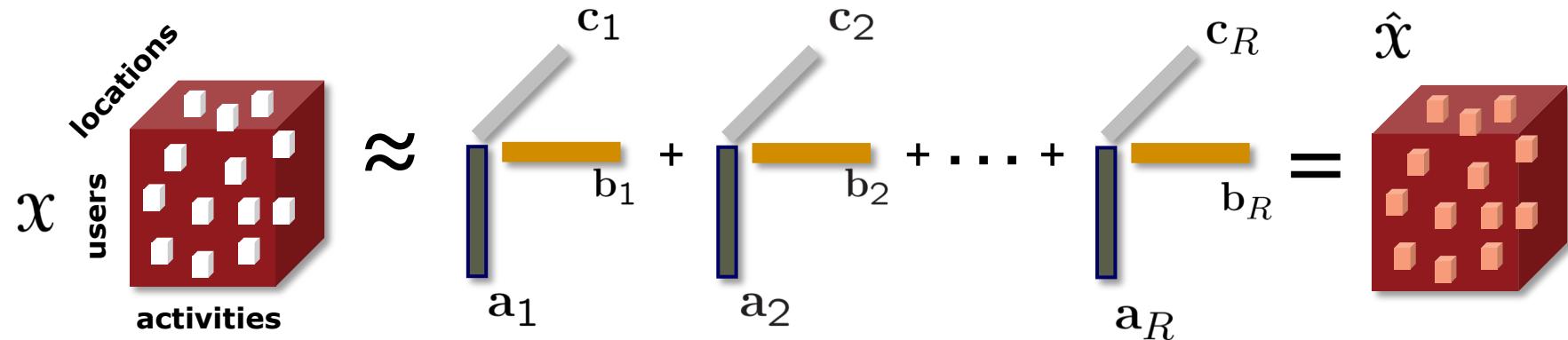


$$\approx \mathbf{U}_3 \begin{matrix} \diagdown \\ \Sigma_3 \end{matrix} \mathbf{V}_3^\top$$



Tensor Factorizations: Missing Data Estimation (a.k.a. Tensor Completion)

Similar to the matrix case, in the case of missing data we can use low-rank tensor approximations to fill in missing entries [Tomasi and Bro, 2005]:



Finding low-rank approximation:

$$\min_{A, B, C} \| \mathcal{W} * (\mathcal{X} - [A, B, C]) \|^2$$

$$w_{ijk} = \begin{cases} 1 & \text{if } x_{ijk} \text{ is known,} \\ 0 & \text{if } x_{ijk} \text{ is missing.} \end{cases}$$

Data reconstruction:

$$\hat{\mathcal{X}} = [A, B, C]$$

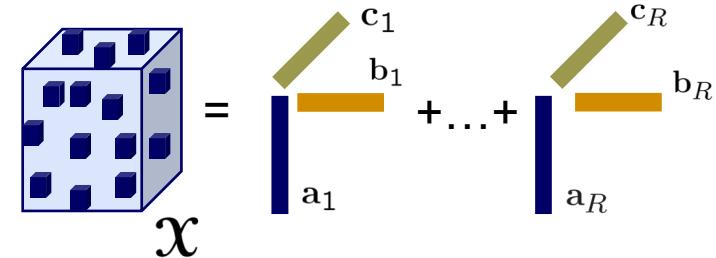
In the presence of missing data: CP-WOPT

[Acar, Dunlavy, Kolda, Mørup, 2011]

CP-WOPT (CP Weighted OPTimization) solves the following weighted optimization problem to fit the CP model to the known data entries in the tensor by ignoring the missing entries:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \| \mathcal{W} * (\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]) \|^2$$

$$w_{ijk} = \begin{cases} 1 & \text{if } x_{ijk} \text{ is known,} \\ 0 & \text{if } x_{ijk} \text{ is missing.} \end{cases}$$



Step1: Define the objective function

$$\begin{aligned} f_{\mathcal{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) &= \| \mathcal{W} * (\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]) \|^2 \\ &= \| \mathbf{y} - \mathbf{z} \|^2 \end{aligned}$$

$$\mathbf{y} = \mathcal{W} * \mathcal{X}$$

$$\mathbf{z} = \mathcal{W} * [\mathbf{A}, \mathbf{B}, \mathbf{C}]$$

Step2: Compute the gradient

$$\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{A}} = 2(\mathbf{z}_{(1)} - \mathbf{y}_{(1)}) (\mathbf{C} \odot \mathbf{B})$$

$$\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{B}} = 2(\mathbf{z}_{(2)} - \mathbf{y}_{(2)}) (\mathbf{C} \odot \mathbf{A})$$

$$\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{C}} = 2(\mathbf{z}_{(3)} - \mathbf{y}_{(3)}) (\mathbf{B} \odot \mathbf{A})$$

→ $\nabla f_{\mathcal{W}} =$
Vectorize and
concatenate
the partials

$$\begin{bmatrix} \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{a}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{a}_R} \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{b}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{b}_R} \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{c}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{c}_R} \end{bmatrix}$$

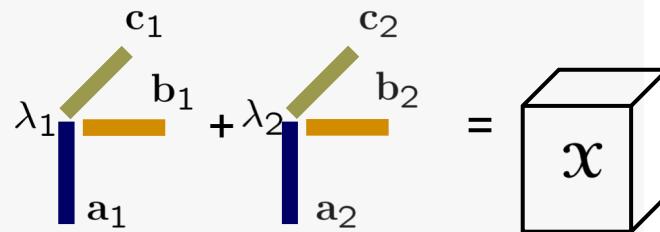


Step3: Pick a first-order optimization method

Tensor Completion Example

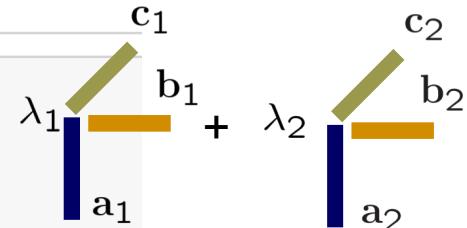
Construct a tensor with an underlying CP model

```
sz = [10 5 3];
R = 2;
for i=1:length(sz)
    A{i}= randn(sz(i),R);
    for r=1:R
        A{i} (:,r)=A{i} (:,r)/norm(A{i} (:,r));
    end
end
lambda = [2 3]';
X = full(ktensor(lambda, A));
Xorig = X.data;
```



Compare the following with ktensor

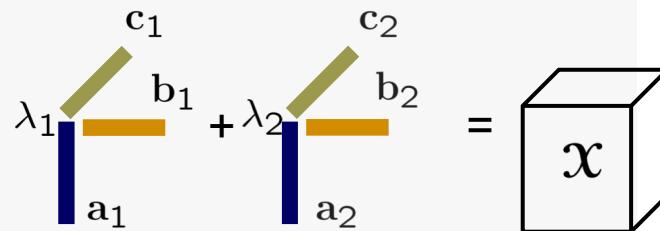
```
Y=zeros(sz);
for r=1:R
    for k=1:size(Y,3)
        Y (:,:,k) = Y (:,:,k) + lambda(r)*A{1} (:,r)*A{2} (:,r)'*A{3} (k,r);
    end
end
```



Tensor Completion Example (cont.)

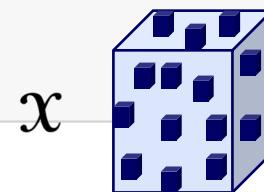
Construct a tensor with an underlying CP model

```
sz = [10 5 3];
R = 2;
for i=1:length(sz)
    A{i} = randn(sz(i),R);
    for r = 1:R
        A{i} (:,r) = A{i} (:,r)/norm(A{i} (:,r));
    end
end
lambda = [2 3]';
X      = full(ktensor(lambda,A));
Xorig = X;
```



Pick some random entries and set them to missing

```
M = round(rand(sz));
X(find(M==0))=NaN;
```



Tensor Completion Example (cont.)

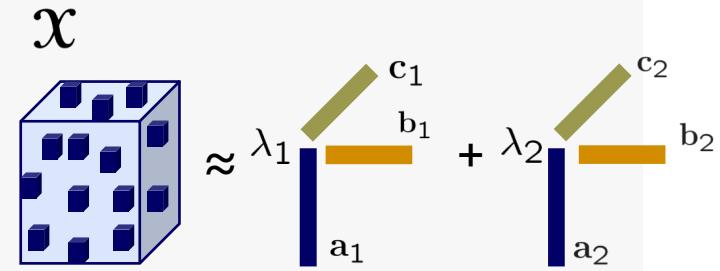
Suppose you are given the incomplete tensor X and you want to find the underlying CP components

```
% Record the places of the missing entries
W = tensor(~isnan(X.data));

% Set those missing entries in the tensor to 0
X(find(W==0))=0;

% parameters for the gradient-based optimization algorithm
options = ncg('defaults');
options.MaxFuncEvals = 10000;
options.MaxIter = 10000; % change this line to change the max number of iterations
options.StopTol = 1e-6;
options.RelFuncTol = 1e-6;
options.DisplayIter = 50;
[Fac, FacInit, output] = cp_wopt(X,W,R,'alg', 'ncg','alg_options',options,'init','nvecs');

% Compare Fac(factors extracted by the model) with A (true factors)
for i=1:3
    corr(Fac.U{i},A{i})
end
```

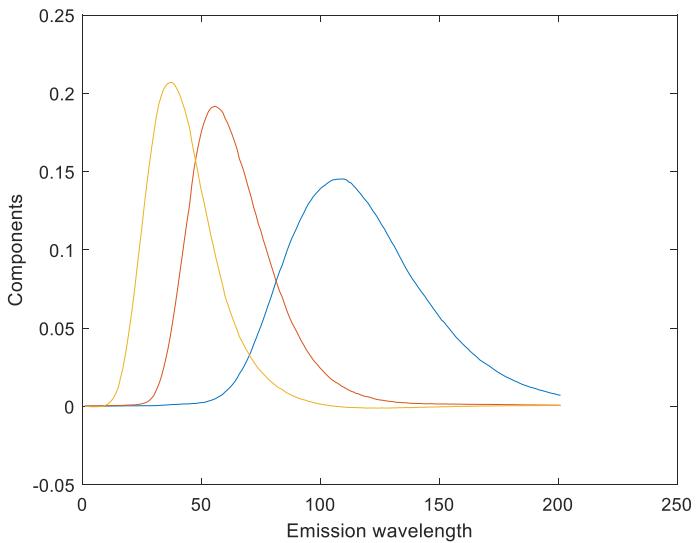
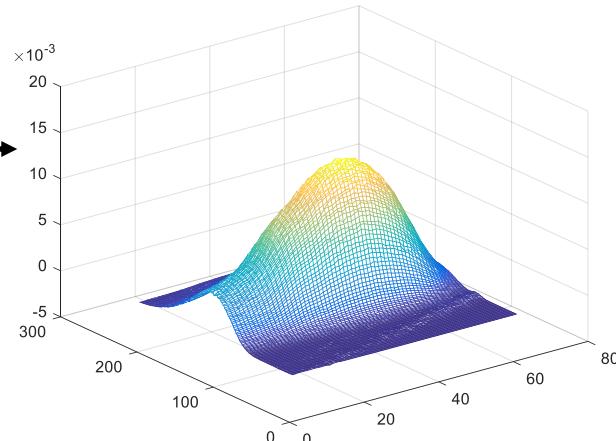
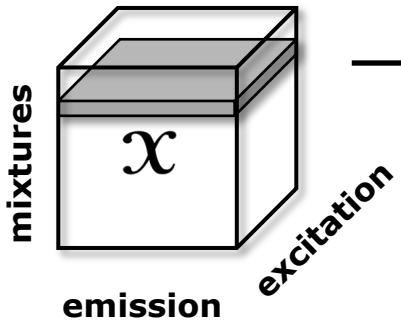


Reconstruct your data from the factor matrices

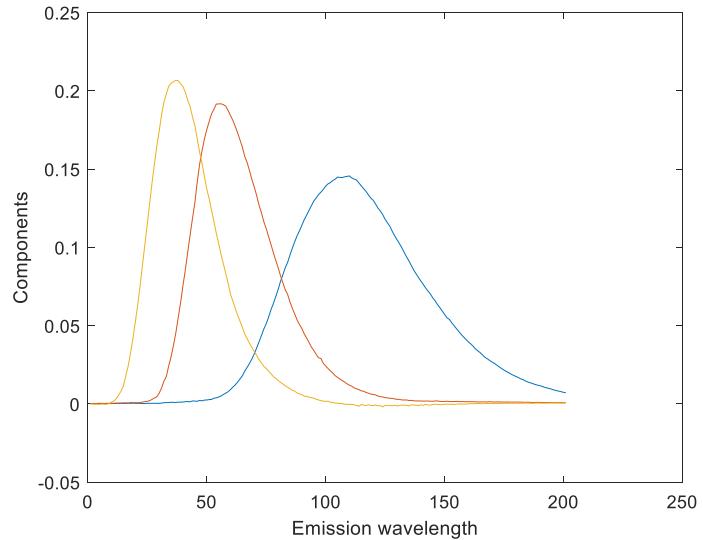
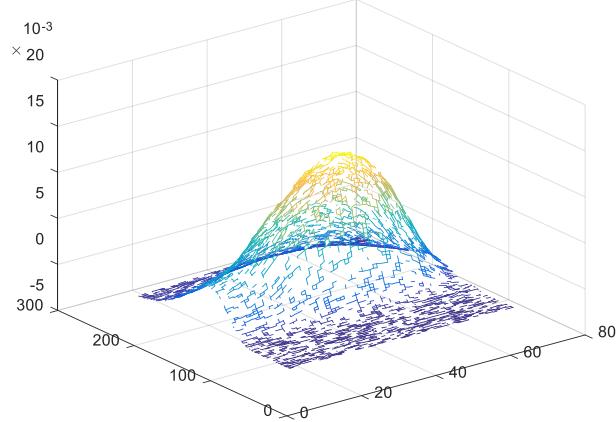
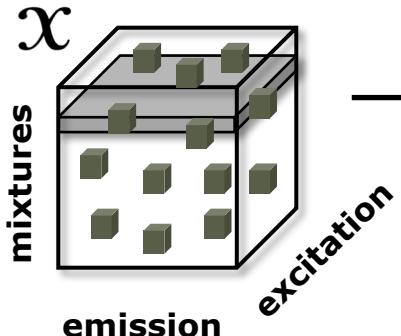
```
Xhat = full(Fac);
% Plot original vs. estimated values for the missing data
plot(Xorig(find(W==0)), Xhat(find(W==0)), 'x');
xlabel('Original'); ylabel('Estimated')
```

$$\hat{X} = [A, B, C]$$

CP-WOPT can accurately capture the factors even with high amounts of missing data!



50% Missing Data

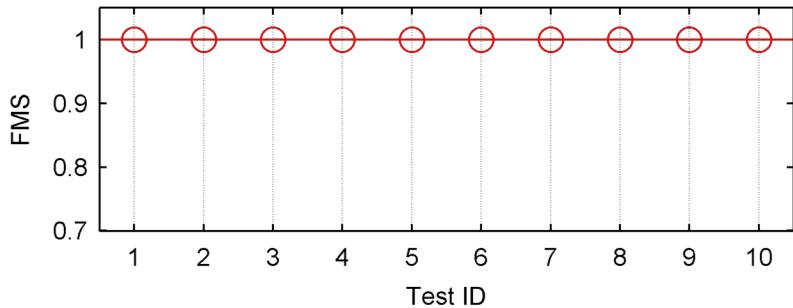
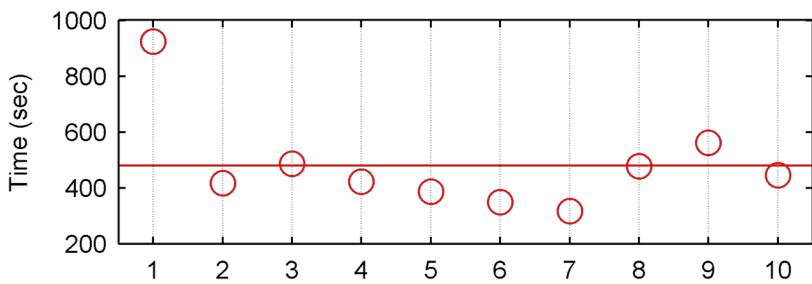


CP-WOPT scales to large problems!

500 x 500 x 500 with 99% missing

entries (i.e., 1.25 million known entries)

**Dense storage = 1GB
Sparse storage = 40MB**

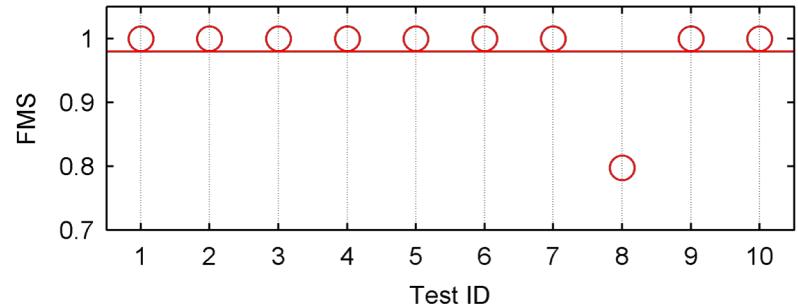
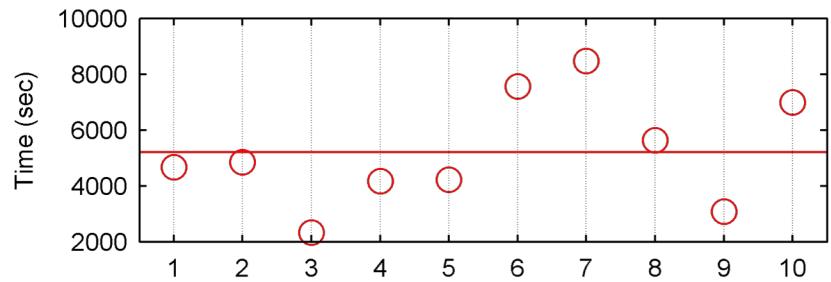


Factor Match Score

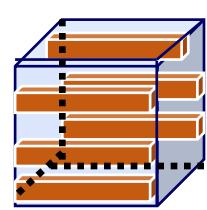
$$FMS = \max_{\sigma \in \Pi(\mathcal{R}, \bar{\mathcal{R}})} \frac{1}{R} \sum_{r=1}^R \left(1 - \frac{|\lambda_r - \bar{\lambda}_{\sigma(r)}|}{\max\{\lambda_r, \bar{\lambda}_{\sigma(r)}\}} \right) \prod_{n=1}^N \left| \bar{a}_r^{(n)\top} \bar{a}_{\sigma(r)}^{(n)} \right|$$

1000 x 1000 x 1000 with 99.5% missing entries (i.e., 5 million known entries)

**Dense storage = 8GB
Sparse storage = 160MB**



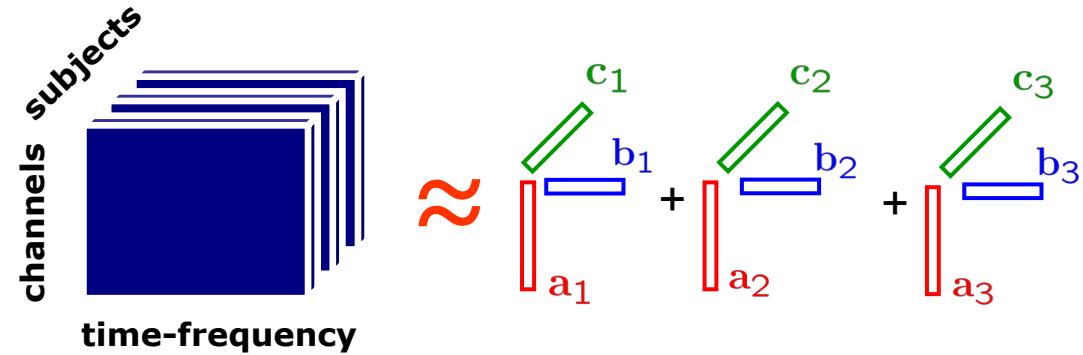
On a Linux workstation with 2 Quad-Core Intel Xeon 3.0 GHz processors and 32 GB RAM

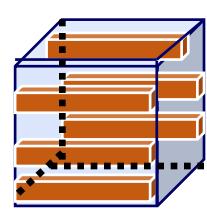


Application₄ (Neuroscience): We can capture underlying factors accurately from EEG data with missing channels

[Acar, Dunlavy, Kolda, Mørup, 2011]

Goal: To differentiate between left and right hand stimulation

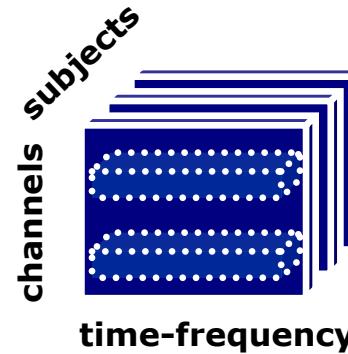




Application₄ (Neuroscience): We can capture underlying factors accurately from EEG data with missing channels

[Acar, Dunlavy, Kolda, Mørup, 2011]

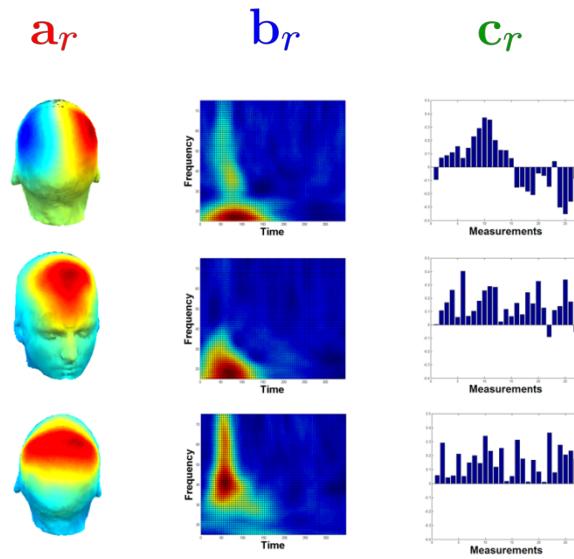
Goal: To differentiate between left and right hand stimulation



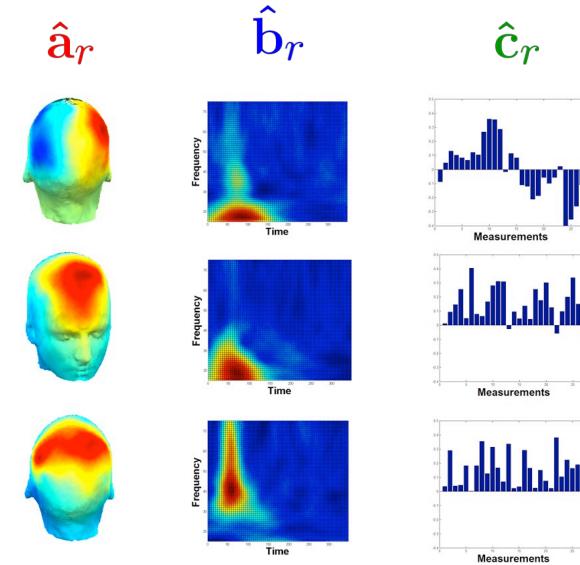
$$\approx c_1 b_1 + c_2 b_2 + c_3 b_3$$

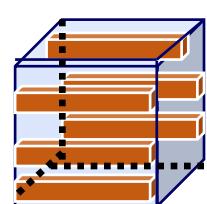
where c_i and b_i are vectors representing components and their corresponding weights.

No Missing Data



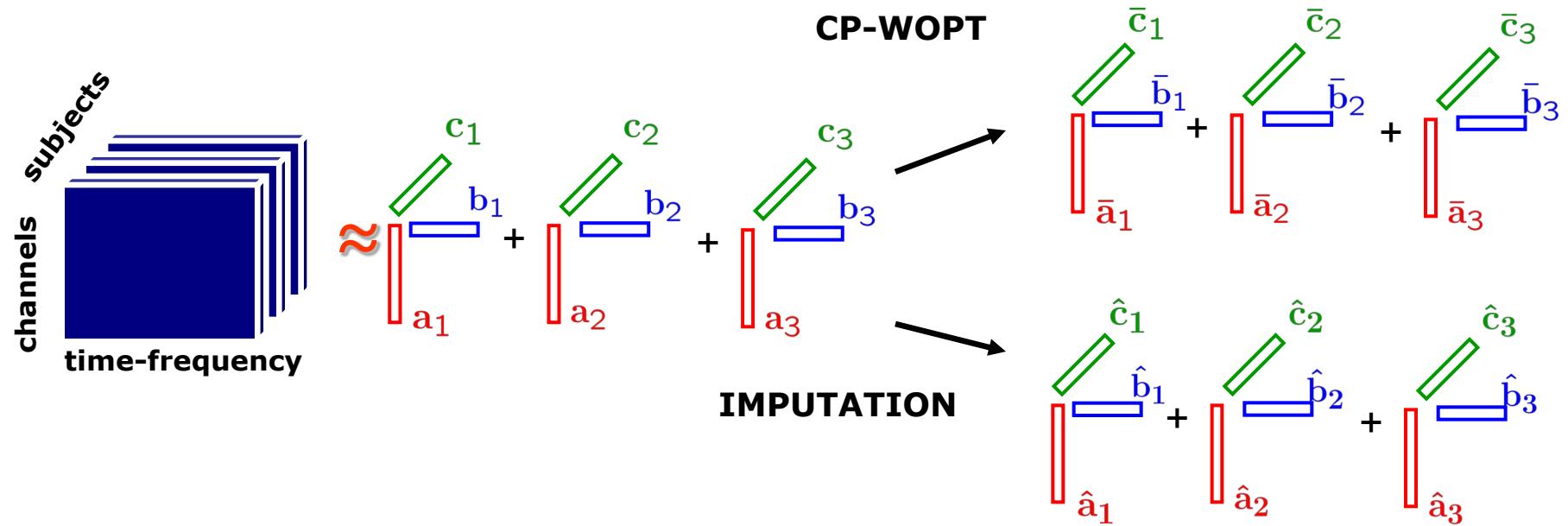
30 Channels/slice Missing

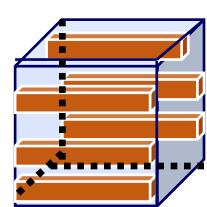




Ignoring missing entries vs. Imputation

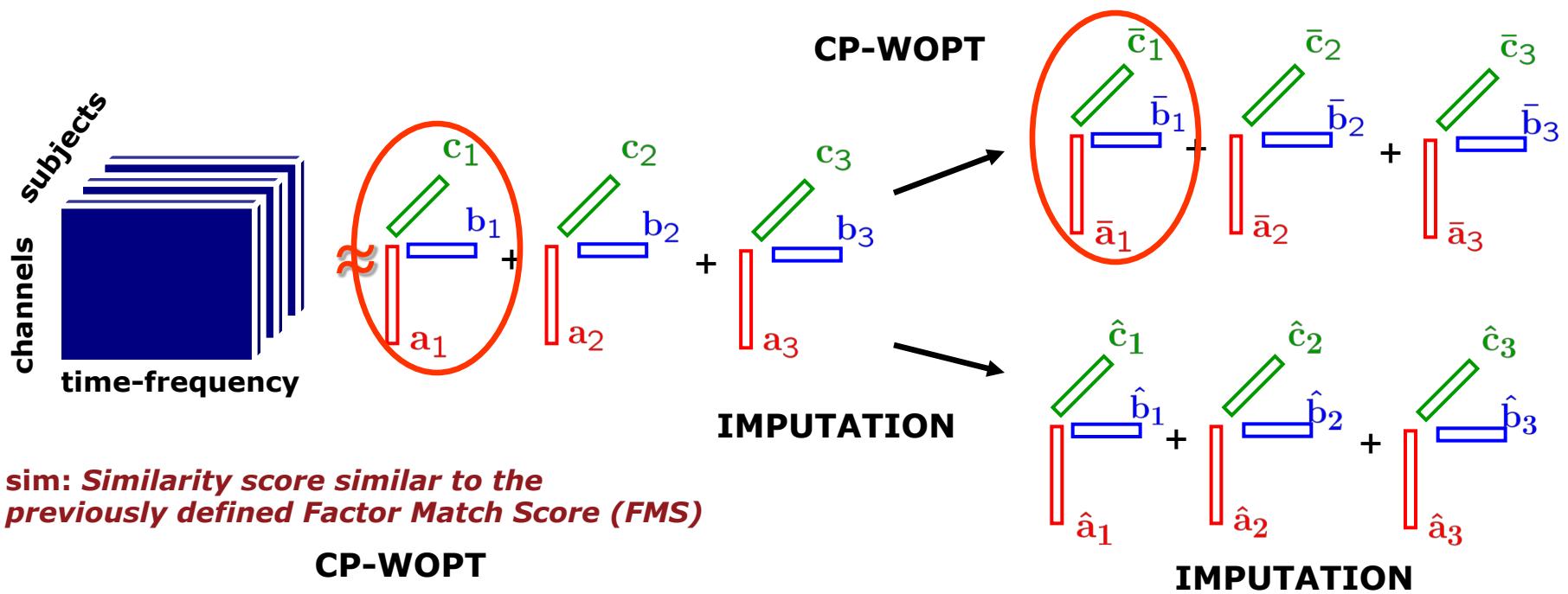
Ignoring missing entries, in particular, when there is large amount of missing data is better than alternatives based on imputation, e.g., imputing missing entries with mean.





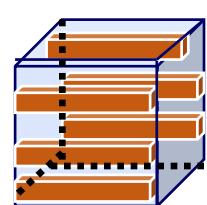
Ignoring missing entries vs. Imputation

Ignoring missing entries, in particular, when there is large amount of missing data is better than alternatives based on imputation, e.g., imputing missing entries with mean.



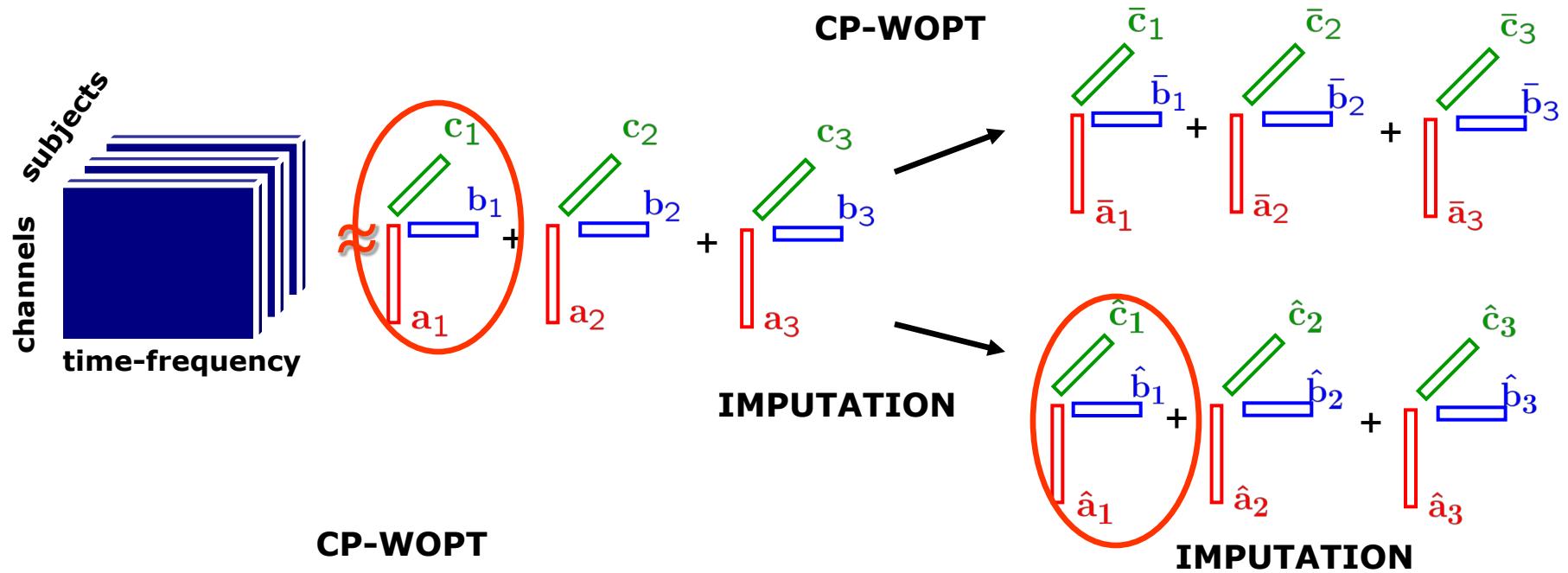
Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9989	0.9995	0.9991
10	0.9869	0.9936	0.9894
20	0.9560	0.9826	0.9697
30	0.9046	0.9604	0.9312
40	0.6192	0.8673	0.7546

Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9970	0.9984	0.9982
10	0.9481	0.9729	0.9752
20	0.9002	0.9475	0.9371
30	0.6435	0.8719	0.8100
40	0.3045	0.7126	0.5699



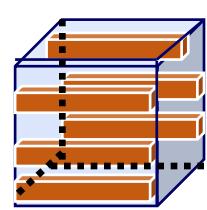
Ignoring missing entries vs. Imputation

Ignoring missing entries, in particular, when there is large amount of missing data is better than alternatives based on imputation, e.g., imputing missing entries with mean.



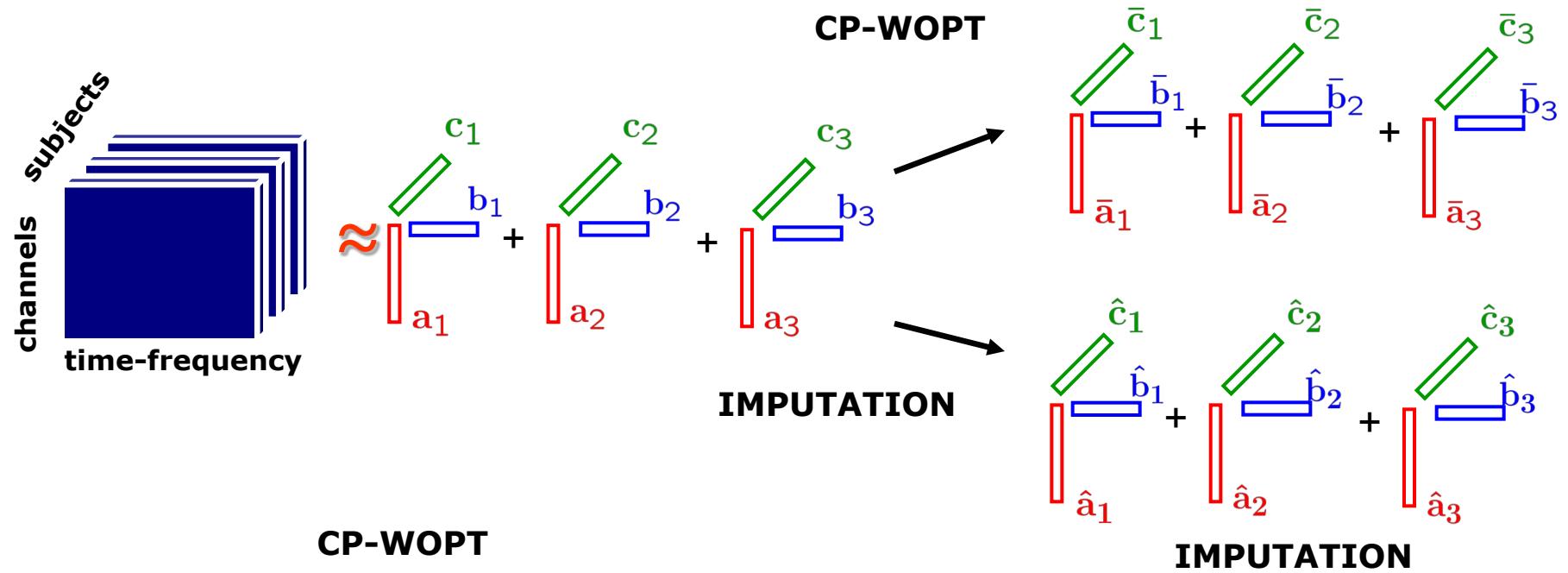
Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9989	0.9995	0.9991
10	0.9869	0.9936	0.9894
20	0.9560	0.9826	0.9697
30	0.9046	0.9604	0.9312
40	0.6192	0.8673	0.7546

Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9970	0.9984	0.9982
10	0.9481	0.9729	0.9752
20	0.9002	0.9475	0.9371
30	0.6435	0.8719	0.8100
40	0.3045	0.7126	0.5699



Ignoring missing entries vs. Imputation

Ignoring missing entries, in particular, when there is large amount of missing data is better than alternatives based on imputation, e.g., imputing missing entries with mean.



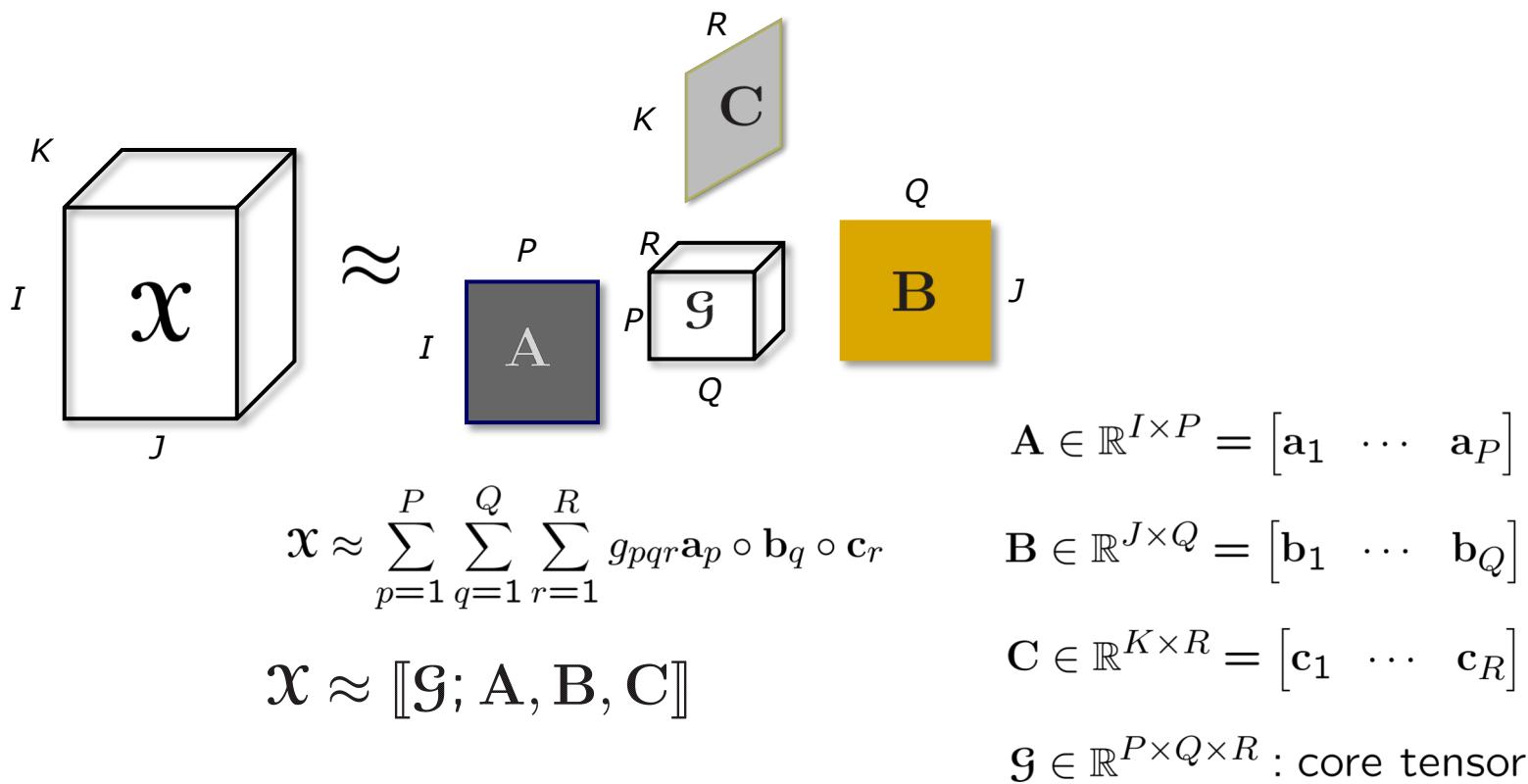
Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9989	0.9995	0.9991
10	0.9869	0.9936	0.9894
20	0.9560	0.9826	0.9697
30	0.9046	0.9604	0.9312
40	0.6192	0.8673	0.7546

Missing Channels	sim($r = 1$)	sim($r = 2$)	sim($r = 3$)
1	0.9970	0.9984	0.9982
10	0.9481	0.9729	0.9752
20	0.9002	0.9475	0.9371
30	0.6435	0.8719	0.8100
40	0.3045	0.7126	0.5699

Another popular tensor factorization model: Tucker

[Tucker, 1963 & 1964 & 1966]

A more flexible tensor factorization model compared to the CP model is the Tucker3 model, which models a third-order tensor using three factor matrices corresponding to each mode as well as a core array.

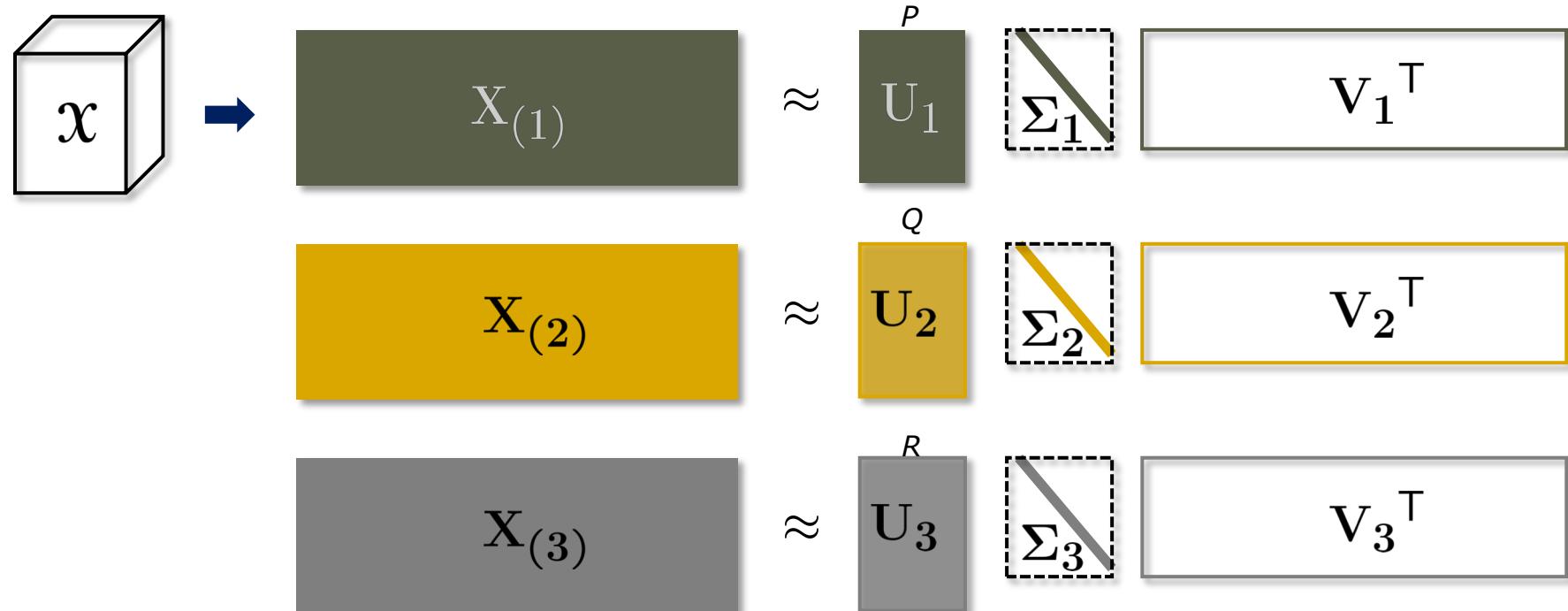
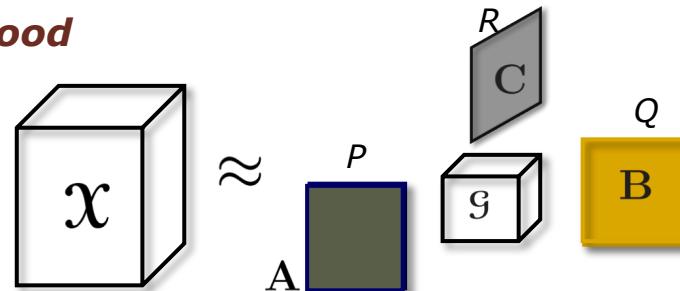


Tucker with orthogonality constraints on A, B, and C corresponds to the best rank-(P,Q,R) approximation of the tensor.

Tucker is considered to be another generalization of SVD to higher-order tensors: Higher-order SVD (HOSVD)

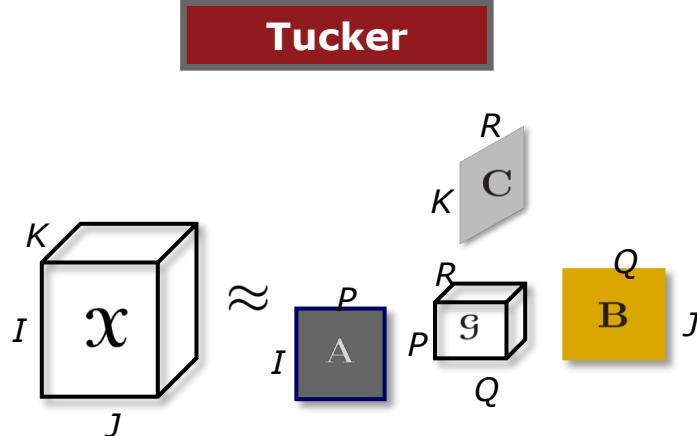
[Tucker, 1966; De Lathauwer, De Moor, Vandewalle, 2000]

Truncation of HOSVD gives a good rank-(P, Q, R) approximation



$$\mathcal{G} = \mathcal{X} \times_1 U_1^\top \times_2 U_2^\top \times_3 U_3^\top$$

Tucker vs. CP



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

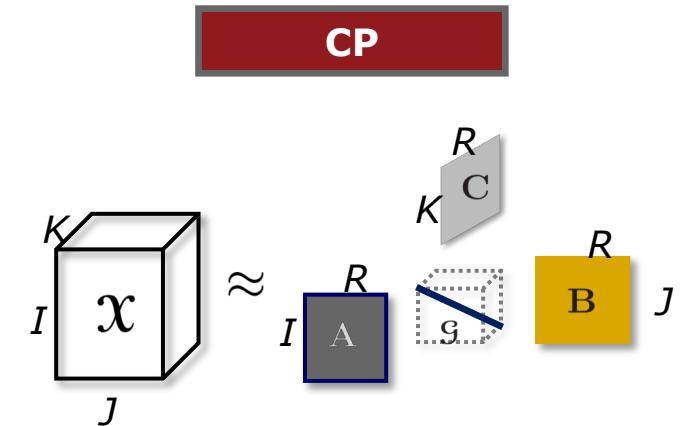
$$\mathcal{X} \approx [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}]$$

More flexible than CP

Different number of components in each mode
The core tensor models the interaction between factors in different modes

Not-unique

$$[\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}] = [\mathcal{G} \times_1 \mathbf{U}; \mathbf{A}\mathbf{U}^{-1}, \mathbf{B}, \mathbf{C}]$$



$$\mathcal{X} \approx \sum_{r=1}^R g_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$$

Restricted model with a specific structure

Same number of components in each mode
Super-diagonal core tensor

Nice uniqueness properties

EXPLORATORY DATA ANALYSIS & COMPRESSION!

INTERPRETABILITY!!!

Application₅ (Face Recognition): TensorFaces

[Vasilescu and Terzopoulos, 2003]

Facial image data tensor with modes (people x viewpoints x illuminations x expressions x pixels) is constructed by arranging the images from 28 people in 5 viewpoints, 4 illuminations and 3 expressions.

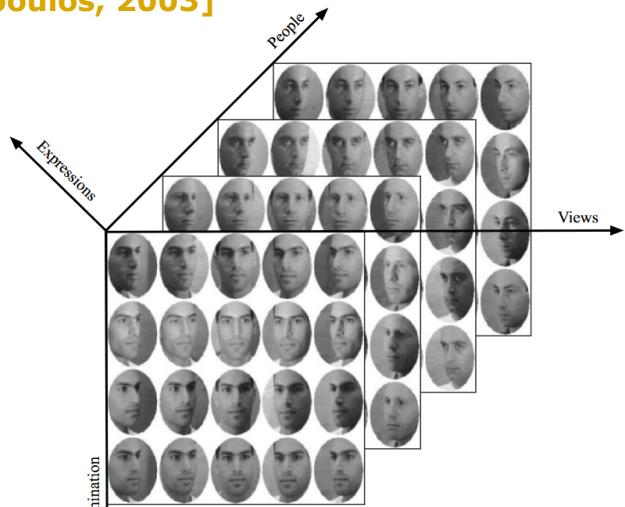
$$\mathcal{D} : 28 \times 5 \times 4 \times 3 \times 7943$$

HOSVD gives the following:

$$\mathcal{D} = \mathcal{Z} \times_1 \mathbf{U}_{\text{people}} \times_2 \mathbf{U}_{\text{views}} \times_3 \mathbf{U}_{\text{illums}} \times_4 \mathbf{U}_{\text{express}} \times_5 \mathbf{U}_{\text{pixels}}$$

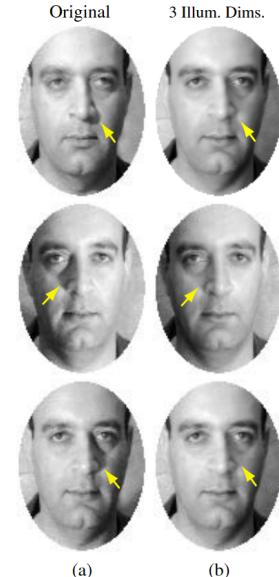
$\mathcal{Z} \times_5 \mathbf{U}_{\text{pixels}}$ **tensorfaces**

$\mathbf{U}_{\text{pixels}}$ **eigenfaces**



Better face recognition performance using tensorfaces compared to eigenfaces

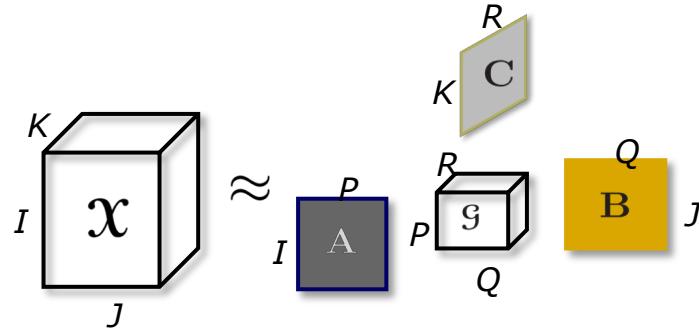
Truncation in the illumination mode:



(a)

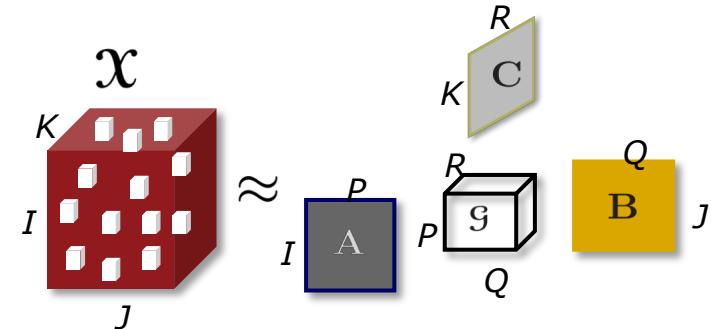
(b)

Similar to CP, Tucker model can be fit to incomplete data using weighted optimization!



$$\mathcal{X} \approx [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}]$$

$$\min_{\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{X} - [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}]\|^2$$



$$\min_{\mathcal{G}, \mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{W} * (\mathcal{X} - [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}])\|^2$$

$$w_{ijk} = \begin{cases} 1 & \text{if } x_{ijk} \text{ is known,} \\ 0 & \text{if } x_{ijk} \text{ is missing.} \end{cases}$$

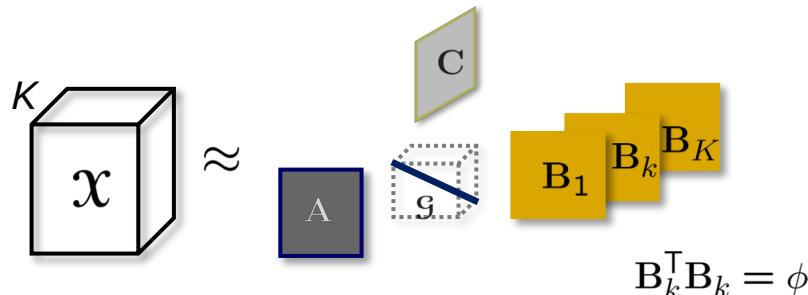
$$\hat{\mathcal{X}} = [\mathcal{G}; \mathbf{A}, \mathbf{B}, \mathbf{C}]$$

There are other tensor factorization models

Mainly for applications with uniqueness & interpretability concerns

PARAFAC2

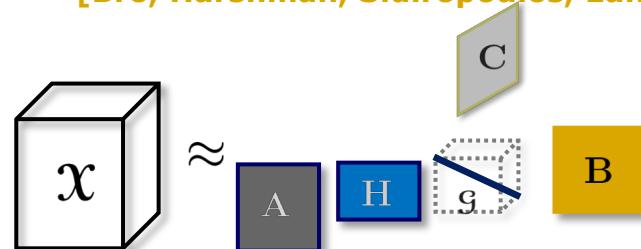
[Harshman, 1972]



Handles differences between factors across different slices. Similar in that sense, there are also Shifted PARAFAC and Convulsive PARAFAC.

PARALIND (Parallel profiles with Linear Dependences)

[Bro, Harshman, Sidiropoulos, Lundy, 2009]



DEDICOM (Decomposition into directional components)

[Harshman, 1978]

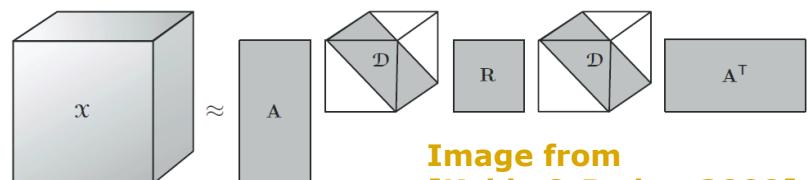


Image from
[Kolda & Bader, 2009]

Mainly for compression & reconstruction applications and higher-order tensors with many modes, e.g., 10th-order tensor

Tensor Train (TT)

And hierarchical tucker (H-Tucker)

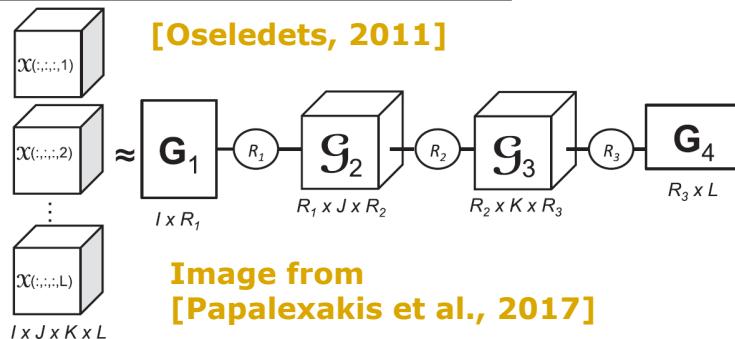


Image from
[Papalexakis et al., 2017]

Sequences of SVDs on various unfoldings:

SVD of $\text{reshape}(\mathcal{X}, [I, JKL]) : G_1, Z_1$

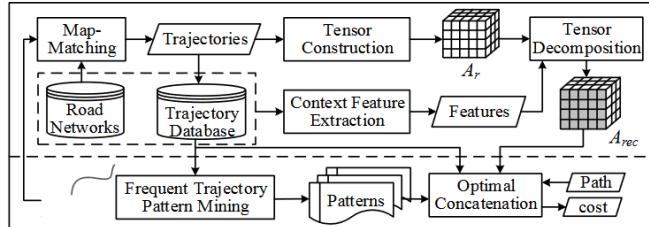
SVD of $\text{reshape}(Z_1, [R_1 J, KL]) : G_2, Z_2$

SVD of $\text{reshape}(Z_2, [R_2 K, L]) : G_3, G_4$

There are many other applications of tensor factorizations

Urban Computing

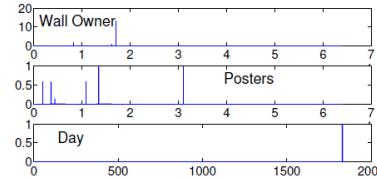
Travel Time Estimation
[Wang et al., KDD, 2014]



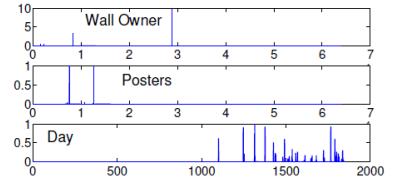
road segments x drivers x time slots

Social Network Analysis

[Papalexakis et al., PKDD, 2012]



(a) FACEBOOK anomaly (Wall owner's birthday)



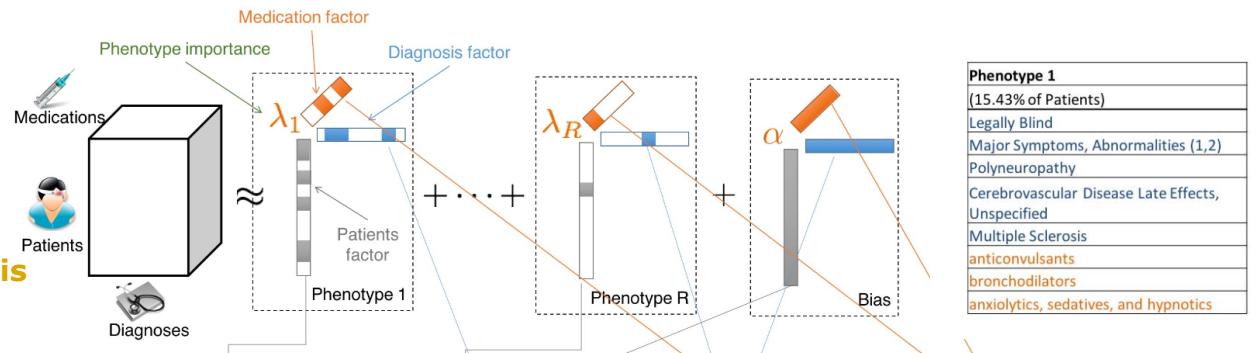
(b) FACEBOOK normal activity

wall owner x poster x day

Electronic Healthcare Records

EHR-based Phenotyping
[Henderson et al., ICHI, 2017]

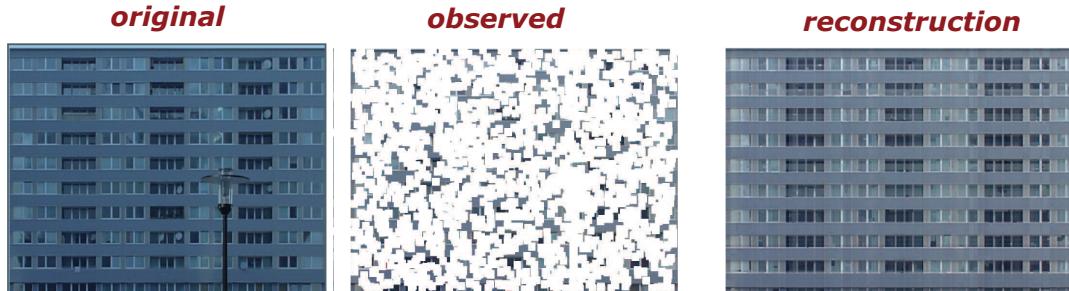
patients x medications x diagnosis



Computer Vision

Image Reconstruction:
[Liu et al., PAMI, 2013]

Tensor: original color image with color channels



Active Research Areas



Tensor Computing for Internet of Things

SCHLOSS DAGSTUHL
Leibniz-Zentrum für Informatik

[Acar, Anandkumar, Mullin,
Rusitschka, Tresp, 2018]

Scalability

Parallel and distributed computations [Smith and Karypis, 2016], [Kaya and Ucar, 2015 & 2016]

Portability to emerging architectures [Phipps et al., 2017]

Sampling [Phan and Cichocki, 2011; Papalexakis et al., 2012; Vervliet and De Lathauwer, 2016; Battaglino et al., 2018]

Different Loss Functions [Yilmaz, Cemgil, Simsekli, 2011; Ermis, Acar, Cemgil, 2015] [Kolda et al., Presented at TRICAP, 2018] [Vandecappelle et al., Presented at TRICAP, 2018]

Deep Learning and Tensor Factorizations

Tensor Train Factorization of a Feed Forward Layer [Novikov, Podoprikhin, Osokin, Vetrov, 2015]

TT-Recurrent Neural Networks [Yang, Krompass, Tresp, 2017]

Software

MATLAB: Nway Toolbox, Tensor Toolbox, TT Toolbox, TensorLab

R: multiway

Sparse: <http://frostd.io/>

Data Sets

Dense: www.models.life.ku.dk

Public data sets for multivariate data analysis

IMPORTANT: all downloadable material listed on these pages - appended by specifics mentioned under the individual headers/chapters - is available for public use. Please note that while great care has been taken, the software, code and data are provided "as is" and that Q&T, LIFE, KU does not accept any responsibility or liability.

Name (click for details)	Description	Located at	Format
NIR	NIR Sugarcane data	University of Copenhagen	Matlab
Tongue data	Three-way data from the work Richard Harshman	University of Copenhagen	Text
JODA data set NEW	NMR, LC-MS and EEM prototypical experimental coupled data sets for JODA	University of Copenhagen	Matlab
RAMAN pork fat NEW	The samples for this study were 16 pork carcasses	University of Copenhagen	Matlab
NIR soil NEW	Soil samples from long-term field experiment in Abisko, northern Sweden	University of Copenhagen	Matlab
AutoChrome NEW	Automatically find PARAFAC2 components of hyphenated chromatographic data	University of Copenhagen	Matlab
Metabolomic data fusion NEW	Biomarker, fluorescence and 1H-NMR data from case/control study on colorectal cancer	University of Copenhagen	Matlab
Olive oil	Oil samples analyzed by HPLC with charged aerosol detector	University of Copenhagen	Matlab
Example tensors	Some example tensors with known problems such as degeneracy, swamps and local minima	University of Copenhagen	Matlab

All tensors:

Name	Non-zeros	Order
Amazon Reviews	1,741,809,018	3
Chicago Crime	5,330,673	4
Delicious	140,126,181	4
Enron Emails	54,202,099	4
Flickr	112,890,310	4
LBNL-Network	1,698,825	5
Matrix Multiplication	M*K*N	3
NELL-1	143,599,552	3
NELL-2	76,879,419	3
NIPS Publications	3,101,609	4
Patents	3,596,640,708	3
Reddit-2015	4,687,474,081	3
Uber Pickups	3,309,490	4
VAST 2015 Mini-Challenge 1	26,021,945	5

Summary

Matrix Factorizations: How to use matrix factorizations to (i) find the underlying factors, (ii) predict missing entries

Matrix factorizations without constraints are not unique.

Constraints need to make sense in terms of the application; otherwise, the model will not reveal what we are looking for.

Tensor Factorizations

Tensor factorization methods such as CP have better uniqueness properties compared to matrix factorizations.

How to use tensor factorizations, in particular, CP, to find the underlying factors, e.g.,

Chemometrics

Neuroscience

Temporal Link Prediction

How to fit a CP model to a higher-order tensor with missing entries

Algorithms: ALS, Gradient-based optimization approaches

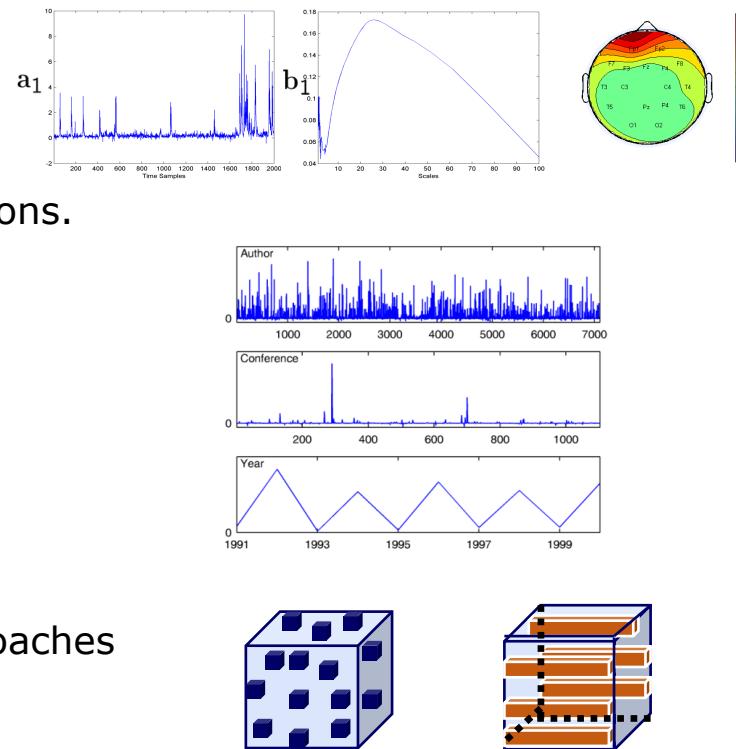
MATLAB functions

cp-opt: Fitting a CP model to a tensor (**from Tensor Toolbox**)

```
[Fac, FaInit, output] = cp_opt(X, R, 'init', 'random', 'alg', 'ncg', 'alg_options', options);
```

cp-wopt: Fitting a CP model to a tensor with missing entries (**from Tensor Toolbox**)

```
[Fac, FaInit, output] = cp_wopt(X, W, R, 'init', 'random', 'alg', 'ncg', 'alg_options', options);
```



Model: Optimization problem

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}] \|^2$$

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \|\mathcal{W} * (\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]) \|^2$$

Tensor Toolbox

Computation of the function value and the gradient

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \|\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}] \|^2$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{a}_1} \\ \vdots \\ \frac{\partial f}{\partial \mathbf{a}_R} \\ \frac{\partial f}{\partial \mathbf{b}_1} \\ \vdots \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{c}_1} \\ \vdots \\ \frac{\partial f}{\partial \mathbf{c}_R} \end{bmatrix} \quad \begin{aligned} \frac{\partial f}{\partial \mathbf{A}} &= -\mathbf{X}_{(1)}(\mathbf{C} \odot \mathbf{B}) + \mathbf{A}(\mathbf{C}^T \mathbf{C} * \mathbf{B}^T \mathbf{B}) \\ \frac{\partial f}{\partial \mathbf{B}} &= -\mathbf{X}_{(2)}(\mathbf{C} \odot \mathbf{A}) + \mathbf{B}(\mathbf{C}^T \mathbf{C} * \mathbf{A}^T \mathbf{A}) \\ \frac{\partial f}{\partial \mathbf{C}} &= -\mathbf{X}_{(3)}(\mathbf{B} \odot \mathbf{A}) + \mathbf{C}(\mathbf{B}^T \mathbf{B} * \mathbf{A}^T \mathbf{A}) \end{aligned}$$

cp_opt

$$\begin{aligned} f_{\mathcal{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) &= \|\mathcal{W} * (\mathcal{X} - [\mathbf{A}, \mathbf{B}, \mathbf{C}]) \|^2 \\ &= \|\mathbf{Y} - \mathbf{Z} \|^2 \end{aligned}$$

$$\nabla f_{\mathcal{W}} = \begin{bmatrix} \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{a}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{a}_R} \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{b}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{b}_R} \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{c}_1} \\ \vdots \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{c}_R} \end{bmatrix} \quad \begin{aligned} \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{A}} &= 2(\mathbf{Z}_{(1)} - \mathbf{Y}_{(1)}) (\mathbf{C} \odot \mathbf{B}) \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{B}} &= 2(\mathbf{Z}_{(2)} - \mathbf{Y}_{(2)}) (\mathbf{C} \odot \mathbf{A}) \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{C}} &= 2(\mathbf{Z}_{(3)} - \mathbf{Y}_{(3)}) (\mathbf{B} \odot \mathbf{A}) \end{aligned}$$

cp_wopt

Poblano Toolbox

Gradient-based optimization algorithms

ncg, lbfgs, ...