# Tensor Factorizations, Data Fusion & Applications

# **Evrim Acar**

Simula Metropolitan Center for Digital Engineering

**Oslo, Norway** 

July 2, 2018, LVA/ICA



# **Matrix Factorizations in Data Mining**





### 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>

#### keywords





### 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>

#### keywords



#### **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!





# **Matrix Factorizations in Data Mining**



**Data sets are often multi-way: Tensor Factorizations** 



Data sets often come from multiple sources: Data Fusion



# Motivation(1): Joint analysis of measurements from multiple platforms has the potential to enhance biomarker discovery

**Neuroscience:** The complexity of the human brain often necessitates the use of **multiple neuroimaging techniques** to better **understand neural activities**. Neuroimaging techniques provide **complementary** information at different scales. **Functional Magnetic** 



Can we jointly analyze these data sets and capture functional/structural patterns that differ between healthy controls and patients suffering from a disorder?

# Motivation(1): Joint analysis of measurements from multiple platforms has the potential to enhance biomarker discovery

**Metabolomics:** The goal is to **detect** a wide range of **metabolites** in biological samples, e.g., blood or urine, and to **discover** the significant metabolic **biomarkers** related to certain conditions such as food intake or various diseases.





# Motivation(1): Joint analysis of measurements from multiple platforms has the potential to enhance biomarker discovery

**Metabolomics:** The goal is to **detect** a wide range of **metabolites** in biological samples, e.g., blood or urine, and to **discover** the significant metabolic **biomarkers** related to certain conditions such as food intake or various diseases.



How can we jointly analyze these coupled data sets and capture the underlying patterns related to various conditions?

# Motivation(2): Data fusion can improve missing data estimation performance

**Recommender Systems:** The goal is to recommend users activities that they may be interested in doing at various locations.



 $\mathbf{Z}$ 

number of points of interest at different locations

### Motivation(2): Data fusion can improve missing data estimation performance

**Recommender Systems:** The goal is to recommend users activities that they may be interested in doing at various locations.



How can we fill in missing entries of the tensor using additional sources of information?



# **Tutorial Outline**

**Matrix Factorizations** 



### **Tensor Factorizations**



### Data Fusion based on Coupled Matrix and Tensor Factorizations



### **Matrix Factorizations**

Matrix Factorizations, e.g., SVD (Singular Value Decomposition), NMF (Nonnegative Matrix Factorization), and Sparse Matrix Factorizations, are commonly used in data mining.





# **Matrix Factorizations: Finding Underlying Structures**

We can use SVD to find the underlying structures in a data set. For instance, using SVD of a *users* by *movies* matrix, we can capture the movie types and user groups.



# **Matrix Factorizations: Finding Underlying Structures**

We can use SVD to find the underlying structures in a data set. For instance, using SVD of a *users* by *movies* matrix, we can capture the movie types and user groups.



#### Who likes what type of movies?



# **Matrix Factorizations: Finding Underlying Structures**



### Matrix Factorizations: Missing Data Estimation (a.k.a. matrix completion)

In the case of incomplete data, matrix factorizations can be used to fill in the missing entries:



Finding low-rank approximation:

$$\min_{\mathbf{A},\mathbf{B}} \left\| \mathbf{W} * (\mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}}) \right\|^2$$

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

**Data reconstruction:** 

$$\hat{\mathbf{X}} = \mathbf{A}\mathbf{B}^{\mathsf{T}}$$

It is possible to recover an unknown lowrank matrix from a nearly minimal set of entries. For more on matrix completion [Candes and Plan, 2010]

# **Matrix Factorizations: Missing Data Estimation**

For instance, we can still factorize the incomplete X and then use those factors to reconstruct the matrix and estimate the missing entries.

	Х		movi	ies				$\hat{\mathbf{X}}$						
					~		$a_1$	<b>b</b> <sub>1</sub>	b a <sub>2</sub>	2 =		4		
X	Aloha	Star Wars	American Pie	Hunger Games	Silver Linings	Maze Runner				Â	$\mathbf{X} = \mathbf{A}$	$\mathrm{B}^{T}$		
User 1	5	2	4	1	3	?	[	3.72	1.27	4.41	1.60	3.78	1.91	
User 2	1	4	1	4	1	4		0.82	3.91	0.74	4.05	1.53	3.91	
User 3	3	1	5	2	5	2		4.06	1.37	4.82	1.72	4.12	2.07	
User 4	1	4	1	?	3	4		1.37	4.06	1.39	4.24	2.08	4.15	
User 5	3	1	4	2	3	2		3.15	1.38	3.71	1.65	3.25	1.90	
User 6	2	?	2	4	2	4		1.77	3.76	1.89	3.97	2.40	3.94	

# **Matrix Factorizations: Missing Data Estimation**

For instance, we can still factorize the incomplete X and then use those factors to reconstruct the matrix and estimate the missing entries.

	Х		mov	ies				$\hat{\mathbf{X}}$					
					~		$a_1$	<b>b</b> <sub>1</sub>	a <sub>2</sub> b	2 =	•	4	
X	Aloha	Star Wars	American Pie	Hunger Games	Silver Linings	Maze Runner				Ŷ	$\mathbf{X} = \mathbf{A}$	$\mathrm{B}^{T}$	
User 1	5	2	4	1	3	2		3.72	1.27	4.41	1.60	3.78	1.91
User 2	1	4	1	4	1	4		0.82	3.91	0.74	4.05	1.53	3.91
User 3	3	1	5	2	5	2		4.06	1.37	4.82	1.72	4.12	2.07
User 4	1	4	1	5	3	4		1.37	4.06	1.39	4.24	2.08	4.15
User 5	3	1	4	2	3	2		3.15	1.38	3.71	1.65	3.25	1.90
User 6	2	4	2	4	2	4		1.77	3.76	1.89	3.97	2.40	3.94

### Matrix factorizations are not unique!



# Uniqueness is an issue

# $\mathbf{X} = \mathbf{A}\mathbf{B}^{\mathsf{T}} = \mathbf{A}\mathbf{M}\mathbf{M}^{-1}\mathbf{B}^{\mathsf{T}} = \bar{\mathbf{A}}\bar{\mathbf{B}}^{\mathsf{T}}$

Constraints are used to deal with the uniqueness problem, e.g., SVD. However, factorizations with constraints may not be meaningful in terms of the application.

**True factors** 





Data matrix

 $X = AB^{\mathsf{T}}$ 



Given X, can we recover the true factors?

### SVD captures...

 $\mathbf{X} \approx \hat{\mathbf{X}} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\mathsf{T}$ 





# Uniqueness is an issue

# $\mathbf{X} = \mathbf{A}\mathbf{B}^\mathsf{T} = \mathbf{A}\mathbf{M}\mathbf{M}^{-1}\mathbf{B}^\mathsf{T} = \bar{\mathbf{A}}\bar{\mathbf{B}}^\mathsf{T}$

Constraints are used to deal with the uniqueness problem, e.g., SVD. However, factorizations with constraints may not be meaningful in terms of the application.

**True factors** 





Data matrix

 $\mathbf{X} = \mathbf{A}\mathbf{B}^{\mathsf{T}}$ 



Given X, can we recover the true factors?

### NMF captures...

 $X\approx \! \hat{X} = WH^{\mathsf{T}}$ 







# Uniqueness is an issue

# $\mathbf{X} = \mathbf{A}\mathbf{B}^\mathsf{T} = \mathbf{A}\mathbf{M}\mathbf{M}^{-1}\mathbf{B}^\mathsf{T} = \bar{\mathbf{A}}\bar{\mathbf{B}}^\mathsf{T}$

Constraints are used to deal with the uniqueness problem, e.g., SVD. However, factorizations with constraints may not be meaningful in terms of the application.

**True factors** 





Data matrix

 $\mathbf{X} = \mathbf{A}\mathbf{B}^{\mathsf{T}}$ 



Given X, can we recover the true factors?

### NMF captures...



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

### **True factors**



$$\mathbf{X}_1 = \mathbf{A} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{B}^\mathsf{T}$$



$$\mathbf{X}_2 = \mathbf{A} \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{B}^\mathsf{T}$$



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

# **Tutorial Outline**

**Matrix Factorizations** 





#### Data Fusion based on Coupled Matrix and Tensor Factorizations



# **Tensors: Terminology**



# **Tensors: Basic Operations**

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



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



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

## **Tensors: Basic Operations (cont.)**

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



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

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

#### **Kronecker product**

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

Khatri-Rao product (Columnwise Kronecker Product)

$$\mathbf{A} \in \mathbb{R}^{I \times R}, \mathbf{B} \in \mathbb{R}^{J \times R}$$
$$\mathbf{A} \odot \mathbf{B} = \begin{bmatrix} \mathbf{a}_1 \otimes \mathbf{b}_1 & \mathbf{a}_2 \otimes \mathbf{b}_2 & \dots & \mathbf{a}_R \otimes \mathbf{b}_R \end{bmatrix} \in \mathbb{R}^{IJ \times R}$$

### **Tensors: Basic Operations (cont.)**

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

$$\langle \mathbf{X}, \mathbf{\mathcal{Y}} \rangle = \sum_{k=1}^{K} \sum_{j=1}^{J} \sum_{i=1}^{I} x_{ijk} y_{ijk} \qquad \langle$$

$$\langle \mathbf{x}, \mathbf{y} \rangle = c$$

#### **Frobenius Norm**

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

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



$$z_{ijk} = x_{ijk} y_{ijk}$$



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

### **True factors**



$$\mathbf{X}_1 = \mathbf{A} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{B}^\mathsf{T}$$



$$\mathbf{X}_2 = \mathbf{A} \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{B}^\mathsf{T}$$



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

# A popular tensor factorization model: CP

### 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)

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:



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

# **CP** is unique up to trivial indeterminacies





# SVD vs. CP

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



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

r=1



# Does it really work?

#### **EXAMPLE:**

#### **True factors**





50

100

150

200

 $\mathbf{X}_2$ 

50

100

150

200

 $\mathbf{X}_1$ 

50

10 20 30 40 50 60

10 20 30 40 50 60

60 60

$$\mathbf{X}_2 = \mathbf{A} \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{B}^\mathsf{T}$$



 $\mathbf{X} \approx \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} 
rbracket$ 



# Algorithms for fitting the CP model



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

$$\mathbf{A} = \mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) ((\mathbf{C} \odot \mathbf{B})^{\mathsf{T}} (\mathbf{C} \odot \mathbf{B}))^{\dagger}$$
$$= \mathbf{X}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^{\mathsf{T}} \mathbf{C} * \mathbf{B}^{\mathsf{T}} \mathbf{B})^{\dagger}$$

Matricized Tensor Times Khatri-Rao Product (MTTKRP)

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

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

while "not converged" do  
Solve for **A** (with fixed **B**, **C**)  

$$\min_{\mathbf{A}} || \mathcal{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] ||^2$$
Solve for **B** (with fixed **A** and **C**)  

$$\min_{\mathbf{B}} || \mathcal{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] ||^2$$
Solve for **C** (with fixed **A** and **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-atonce 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.

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



**Step1: Define the objective function** 

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



### 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**  $\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$   $\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \ddots & \vdots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix}$ 



#### **First-order Necessary Conditions**

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

#### **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:



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

#### **Taylor's approximation:**

Suppose that *f* is twice continuously differentiable, then,

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

**Steepest descent:** 

 $p_k = -\nabla f_k$ 

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

**Newton's direction:** 

**Nonlinear Conjugate Gradient Methods:** 

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

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

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

# Application<sub>1</sub> (Neuroscience): CP has proved useful in epileptic seizure localization [Acar et al., 2007; De Vos et al., 2007]







40 50 60 70 80 90 100 Scales

# Application<sub>2</sub> (Recommender Systems): CP can capture temporal patterns useful for link prediction

[Dunlavy, Kolda, Acar, 2011]

**Temporal Link Prediction** 



## Application<sub>3</sub> (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/



## Application<sub>3</sub> (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.





# simula

# 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.RelFuncTol = 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 (lbfgs)

#### **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.



### **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:

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

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

**Data reconstruction:** 

$$\hat{\mathfrak{X}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} 
brace$$

# 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:

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \| \mathbf{\mathcal{W}}_{*}(\mathbf{\mathcal{X}} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket) \|^{2}$$
$$w_{ijk} = \begin{cases} 1 & \text{if } x_{ijk} \text{ is known,} \\ 0 & \text{if } x_{ijk} \text{ is missing.} \end{cases}$$



### **Step1: Define the objective function**

$$f_{\mathcal{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \| \mathcal{W} * (\mathcal{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket) \|^2$$
$$= \| \mathcal{Y} - \mathcal{Z} \|^2$$

$$\begin{split} \mathcal{Y} &= \mathcal{W} * \mathcal{X} \\ \mathcal{Z} &= \mathcal{W} * \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket \end{split}$$

Step2: Compute the gradient  $\frac{\partial f_{W}}{\partial A} = 2(Z_{(1)} - Y_{(1)})(C \odot B)$   $\frac{\partial f_{W}}{\partial B} = 2(Z_{(2)} - Y_{(2)})(C \odot A)$   $\frac{\partial f_{W}}{\partial B} = 2(Z_{(3)} - Y_{(3)})(B \odot A)$   $\int Vectorize and concatenate the partials$   $\int \frac{\partial f_{W}}{\partial c_{1}} = \frac{\partial f_{W}}{\partial c_{R}}$ 

Step3: Pick a first-order optimization method

# **Tensor Completion Example**

### Construct a tensor with an underlying CP model



#### Compare the following with ktensor



# **Tensor Completion Example (cont.)**

#### Construct a tensor with an underlying CP model



#### Pick some random entries and set them to missing





# **Tensor Completion Example (cont.)**

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



#### 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)), '*');
xlabel('Original'); ylabel('Estimated')
\hat{X} = \llbracket A, B, C \rrbracket
```

# **CP-WOPT** can accurately capture the factors even with high amounts of 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

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













# Application<sub>4</sub> (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





# Application<sub>4</sub> (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





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



### **CP-WOPT**

Missing	sim(r = 1)	sim(r = 2)	sim(r = 3)
Channels			
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

#### IMPUTATION

Missing	sim(r=1)	sim(r = 2)	sim(r = 3)
Channels			
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, 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	sim(r=1)	sim(r=2)	sim(r = 3)
Channels			
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, in particular, when there is large amount of missing data is better than alternatives based on imputation, e.g., imputing missing entries with mean.



#### **CP-WOPT**

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

#### **IMPUTATION**

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.

$$\begin{array}{c} \overset{\kappa}{\overbrace{}} & \overset{\kappa}{\overbrace{}} & \overset{\kappa}{\overbrace{}} & \overset{Q}{\overbrace{}} \\ I & \overset{\chi}{\overbrace{}} & \overset{\varphi}{\underset{I}} & \overset{P}{\overbrace{}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{Q}{\underset{Q}} \\ I & \overset{P}{\overbrace{}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{Q}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{Q}{\underset{Q}} \\ I & \overset{R}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}} \\ I & \overset{P}{\underset{Q}} & \overset{P}{\underset{Q}}$$

# 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]



 $\mathbf{9} = \mathbf{X} \times_1 \mathbf{U}_1^{\mathsf{T}} \times_2 \mathbf{U}_2^{\mathsf{T}} \times_3 \mathbf{U}_3^{\mathsf{T}}$ 

# Tucker vs. CP



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

### More flexible than CP

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

#### Not-unique

$$\llbracket \mathbf{9}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket = \llbracket \mathbf{9} \times_1 \mathbf{U}; \mathbf{A}\mathbf{U}^{-1}, \mathbf{B}, \mathbf{C} \rrbracket$$

# EXPLORATORY DATA ANALYSIS & COMPRESSION!



### Restricted model with a specific structure

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

#### Nice uniqueness properties

INTERPRETABILITY!!!

# **Application<sub>5</sub>** (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.

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

HOSVD gives the following:

```
\mathcal{D} = \mathcal{Z} \times_{1} U_{\text{people}} \times_{2} U_{\text{views}} \times_{3} U_{\text{illums}} \times_{4} U_{\text{expres}} \times_{5} U_{\text{pixels}}
```



Better face recognition performance using tensorfaces compared to eigenfaces

Truncation in the illumination mode:



Views

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





 $\mathfrak{X} \approx \llbracket \mathfrak{G}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ 

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

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

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

 $\hat{\mathbf{X}} = \llbracket \mathbf{G}; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ 

# simula

## There are other tensor factorization models

### Mainly for applications with uniqueness & interpretability concerns



Handles differences between factors across different slices. Similar in that sense, there are also Shifted PARAFAC and Convolutive PARAFAC. PARALIND (Parallel profiles with Linear Dependences)

[Bro, Harshman, Sidiropoulos, Lundy, 2009]



DEDICOM (Decomposition into directional components [Harshman, 1978]



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



And hierarchical tucker (H-Tucker)

Sequences of SVDs on various unfoldings:

SVD of reshape( $\mathfrak{X}, [I, JKL]$ ) :  $\mathbf{G}_1, \mathbf{Z}_1$ 

- SVD of reshape( $\mathbf{Z}_1$ , [ $R_1J$ , KL]) :  $\mathbf{G}_2$ ,  $\mathbf{Z}_2$
- SVD of reshape $(\mathbf{Z}_2, [R_2K, L])$  :  $\mathfrak{G}_3, \mathbf{G}_4$

## 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



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

#### wall owner x poster x day





(b) FACEBOOK normal activity



#### **Computer Vision**

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

Tensor: original color image with color channels

#### original



#### reconstruction



## **Active Research Areas**

#### Manifesto from Dagstuhl Perspectives Workshop 16152

Tensor Computing for Internet of Things

S C H L O S S D A G S T U H L 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://frostt.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, LITE, 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 NEV	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

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

#### All tensors:

## 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, FacInit, 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, FacInit, output] = cp\_wopt(X, W, R, 'init', 'random', 'alg','ncg', 'alg\_options', options);







Model: Optimization problem

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C}} \| \boldsymbol{\mathfrak{X}} - [\![\mathbf{A},\mathbf{B},\mathbf{C}]\!] \|^2 \qquad \min_{\mathbf{A},\mathbf{B},\mathbf{C}} \| \boldsymbol{\mathscr{W}} * (\boldsymbol{\mathfrak{X}} - [\![\mathbf{A},\mathbf{B},\mathbf{C}]\!]) \|^2$$

Tensor Toolbox Computation of the function value and the gradient

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \| \mathbf{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \|^2 \qquad f_{\mathcal{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \| \mathcal{W} * (\mathbf{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]) \|^2 \\ = \| \mathcal{Y} - \mathcal{Z} \|^2 \\ = \| \mathcal{Y} - \mathcal{Z} \|^2 \\ = \| \mathcal{Y} - \mathcal{Z} \|^2 \\ \frac{\partial f}{\partial \mathbf{a}_R} \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{c}_1} \\ \frac{\partial f}{\partial \mathbf{c}_R} \\ \frac{\partial f}{\partial \mathbf{c}_R} \\ \frac{\partial f}{\partial \mathbf{c}_R} = -\mathbf{X}_{(3)}(\mathbf{C} \odot \mathbf{A}) + \mathbf{B}(\mathbf{C}^{\mathsf{T}}\mathbf{C} * \mathbf{A}^{\mathsf{T}}\mathbf{A}) \\ \mathbf{Cp\_opt} \qquad \nabla f_{\mathcal{W}} = \begin{bmatrix} \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{a}_R} \\ \frac{\partial f_{\mathcal{W}}}{\partial \mathbf{b}_R} \\ \frac{\partial f_{$$

Poblano Toolbox Gradient-based optimization algorithms

ncg, lbfgs, ...

## **Tutorial Outline**

**Matrix Factorizations** 



#### **Tensor Factorizations**





# **Data Fusion based on Coupled Factorizations**



$$\min_{\mathbf{U},\mathbf{V},\mathbf{W}} \left\| \mathbf{X} - \mathbf{U}\mathbf{V}^{\mathsf{T}} \right\|^{2} + \left\| \mathbf{Y} - \mathbf{U}\mathbf{W}^{\mathsf{T}} \right\|^{2}$$

**Psychometrics:** Simultaneous factorization of Gramian matrices [Levin, 1966]; Simultaneous Component Analysis [Kiers and Ten Berge, 1994]

**Chemometrics:** Principal Component Analysis of multiple matrices [Westerhuis et al., 1998]; Augmented Multivariate Curve Resolution [De Juan and Tauler, 2000]

**Bioinformatics:** Comparing genome-scale expression data from multiple organisms [Alter et al., 2003; Ponnapalli et al., 2011]; Clustering microarray data [Badea, 2007]; Simultaneous Component Analysis with rotation to common and distinct components [Van Deun et al., 2012]; Decomposing multiple matrices into terms explaining joint and individual variation [Lock et al., 2013]

**Signal Processing:** Joint diagonalization of multiple matrices [Yeredor, 2002; Ziehe et al., 2004]; Audio source separation [Yoo et al., 2010]; Joint Independent Component Analysis [Calhoun et al., 2006]; Soft nonnegative matrix co-factorization [Seichepine et al., 2013]

**Data Mining:** Collective matrix factorization [Singh and Gordon, 2008]; Clustering multi-type relational data [Long et al., 2006]; Social recommendation [Ma et al., 2008]; Coupled matrix factorization with sparse factors [Van Deun et al., 2011; Acar et al., 2012]; Nonnegative shared subspace learning [Gupta et al., 2010]; Bayesian interbattery factor analysis [Klami et al., 2013]



# **Data Fusion based on Coupled Factorizations**



$$\min_{\mathbf{U},\mathbf{V},\mathbf{W}} \left\| \mathbf{X} - \mathbf{U}\mathbf{V}^{\mathsf{T}} \right\|^{2} + \left\| \mathbf{Y} - \mathbf{U}\mathbf{W}^{\mathsf{T}} \right\|^{2}$$

psychometrics, chemometrics, bioinformatics, signal processing, data mining, ...

*Cannot handle joint analysis of matrices and higher-order tensors!* 



Psychometrics: Linked-mode PARAFAC [Harshman and Lundy, 1984]

Chemometrics: Multi-way Multi-block component models [Smilde et al., 2000]

Bioinformatics: Coupled analysis of in vitro and histology tissue samples [Acar et al., 2012]

**Signal Processing:** Joint analysis of a covariance matrix and a cumulant tensor [De Lathauwer and Vandewalle, 2004; Comon, 2004]; Generalized Coupled Tensor Factorizations [Yilmaz et al., 2011]; Structured Data Fusion [Sorber et al., 2015]

**Data Mining:** Multi-way Clustering [Banerjee et al., 2007]; Community detection [Lin et al., 2009]; Missing value estimation [Zheng et al., 2010]; Link prediction [Ermis et al., 2012]; Scalable CMTF approaches (sampling-based [Papalexakis et al., 2014], distributed stochastic gradient running on MAPREDUCE [Beutel et al., 2014], distributed ALS running on MAPREDUCE [Jeon et al., 2016])

# **Coupled Matrix Factorizations (CMF)**



Use a first-order optimization method to fit the CMF model

# simula

## **Coupled Matrix Factorizations: Example**



$$\mathbf{X} = \mathbf{A}\mathbf{B}^{\mathsf{T}} + \eta \mathbf{N}_{\mathbf{X}}$$
$$\mathbf{Y} = \mathbf{A}\mathbf{C}^{\mathsf{T}} + \eta \mathbf{N}_{\mathbf{Y}}$$

**Construct coupled matrices** 

```
sz = [30 20 10];
modes = {[1 2], [1 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
X{1} = A{1}*A{2}';
X{2} = A{1}*A{3}';
% add noise
% for i=1:2
옿
     N{i} = randn(size(X{i}));
     X{i} = X{i} + 0.1* N{i}/norm(N{i}, 'fro')*norm(X{i}, 'fro');
÷.
% end
P = length(X);
for p=1:P
    Z.object{p} = tensor(X{p});
    Z.object{p} = Z.object{p}/norm(Z.object{p});
end
Z.modes = modes;
Z.size = sz;
```

# simula
## **Coupled Matrix Factorizations: Example (cont.)**



#### Suppose you are given coupled matrices and want to jointly factorize them



#### Compare the reconstructed matrices with the originals

```
Xhat{1} = Fac.U{1}*Fac.U{2}';
Xhat{2} = Fac.U{1}*Fac.U{3}';
norm(Z.object{1}.data-Xhat{1},'fro')
norm(Z.object{2}.data-Xhat{2},'fro')
```

## CMF is not unique!

Compare the factor matrices corresponding to the same function value

```
options.Display ='final';
[Fac1, FacInit1, out1] = cmtf_opt(Z,R,'alg', 'ncg','alg_options',options,'init','random');
[Fac2, FacInit2, out2] = cmtf_opt(Z,R,'alg', 'ncg','alg_options',options,'init','random');
for i=1:3
    corr(Fac1.U{i},Fac2.U{i})
    pause
end
```





## CMF extends to joint analysis of incomplete matrices!



**Objective function:**  $f_{W}(A, B, C) = \underbrace{\left\| W * (X - AB^{\mathsf{T}}) \right\|^{2}}_{f_{W_{1}}} + \underbrace{\left\| Y - AC^{\mathsf{T}} \right\|^{2}}_{f_{W_{2}}}$ 

**Compute the gradient:** 

$$\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{A}} = \frac{\partial f_{\mathbf{W}_1}}{\partial \mathbf{A}} + \frac{\partial f_{\mathbf{W}_2}}{\partial \mathbf{A}}$$
$$\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{B}} = \frac{\partial f_{\mathbf{W}_1}}{\partial \mathbf{B}}$$
$$\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{C}} = \frac{\partial f_{\mathbf{W}_2}}{\partial \mathbf{C}}$$

Use a first-order optimization method to fit the CMF model

## **Coupled Matrix Factorizations: Missing Data Estimation**

Generate factor matrices



## **Coupled Matrix Factorizations: Missing Data Estimation (cont.)**



## CMF can jointly analyze multiple matrices coupled in different modes!



$$\min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D}} \left\| \mathbf{W} * (\mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}}) \right\|^{2} + \left\| \mathbf{Y} - \mathbf{A}\mathbf{C}^{\mathsf{T}} \right\|^{2} + \left\| \mathbf{Z} - \mathbf{D}\mathbf{B}^{\mathsf{T}} \right\|^{2}$$



## **Data Fusion based on Coupled Matrix and Tensor Factorizations**

Joint analysis of heterogeneous data from multiple sources can be formulated as a coupled matrix and tensor factorization (CMTF) problem. In CMTF, higher-order tensors and matrices are simultaneously factorized by fitting a CP model to higher-order tensors and factorizing matrices in a coupled manner.



## All-at-once Optimization for CMTF

[Acar, Dunlavy, Kolda, 2011]

 $\partial f$ 

 $\overline{\partial \mathbf{a}_R} \ \overline{\partial f} \ \overline{\partial \mathbf{b}_1}$ 

 $rac{\partial f}{\partial \mathbf{b}_R}$ 

 $\frac{\partial f}{\partial \mathbf{c}_1}$ 

 $\vdots$  $\partial f$ 

 $\overline{\partial \mathbf{c}_R}$ 

 $\begin{array}{c} \frac{\partial f}{\partial \mathbf{d}_1} \\ \vdots \\ \partial f \end{array}$ 

 $\partial \check{\mathrm{d}}_R$ 

 $\nabla f \equiv$ 

**CMTF-OPT** is a gradient–based optimization approach for joint factorization of coupled matrices and higher-order tensors.

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D}} \| \boldsymbol{\mathfrak{X}} - [\![\mathbf{A},\mathbf{B},\mathbf{C}]\!] \|^2 + \| \mathbf{Y} - \mathbf{A}\mathbf{D}^\mathsf{T} \|^2$$

## **Step 1: Define the objective function**

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \frac{1}{2} \| \mathbf{X} - [ [\mathbf{A}, \mathbf{B}, \mathbf{C} ] ] \|^2 + \frac{1}{2} \| \mathbf{Y} - \mathbf{A} \mathbf{D}^{\mathsf{T}} \|^2 \int_{\mathbb{R}^3 \mathbf{a}_1}^{\frac{\partial f}{\partial \mathbf{a}_1}}$$

## Step 2: Compute the gradient

## Step 3: Pick a first-order optimization method



$$\mathbf{Y} pprox \mathbf{A} \mathbf{D}^{\mathsf{T}} \ \mathbf{\mathfrak{X}} pprox \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} 
rbracket$$



## **Coupled Matrix and Tensor Factorizations: Example**

## **Coupled Matrix and Tensor Factorizations: Example (cont.)**





## **CMTF** inherits uniqueness from CP!

[Sorensen and De Lathauwer, 2015]



#### Compare the factor matrices corresponding to the same function value

```
options.Display ='final';
[Fac1, FacInit1, out1] = cmtf_opt(Z,R,'alg', 'ncg','alg_options',options,'init','random');
[Fac2, FacInit2, out2] = cmtf_opt(Z,R,'alg', 'ncg','alg_options',options,'init','random');
for i=1:4
    corr(Fac1.U{i},Fac2.U{i})
end
```

## **CMTF** easily extends to incomplete data sets

We fit the model only to the known data entries and ignore the missing entries (in higher-order tensors and/or matrices).

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D}} \| \mathbf{\mathcal{W}}_{*}(\mathbf{\mathcal{X}} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket) \|^{2} + \| \mathbf{Y} - \mathbf{A}\mathbf{D}^{\mathsf{T}} \|$$

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



 $\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{a}_1}$ 

 $rac{\partial f_{\mathbf{W}}}{\partial \mathbf{a}_R}$ 

 $rac{\partial f_{\mathbf{W}}}{\partial \mathbf{b}_1}$ 

 $\partial f_{\mathcal{W}}$ 

 $\frac{\partial \mathbf{b}_R}{\partial \mathbf{b}_R}$ 

 $\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{c}_1}$ 

 $\frac{\partial f_{\mathbf{W}}}{\partial \mathbf{c}_{B}}$ 

 $rac{\partial f_{\mathbf{W}}}{\partial \mathbf{d}_1}$ 

 $rac{\partial f_{\mathbf{W}}}{\partial \mathbf{d}_R}$ 

. . .

 $\nabla f_{\mathbf{w}} =$ 

Vectorize and

concatenate

the partials

2

### **Our objective:**

$$f_{\mathcal{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \frac{1}{2} \| \mathcal{W}_{*}(\mathcal{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!]) \|^{2} + \frac{1}{2} \| \mathbf{Y} - \mathbf{A}\mathbf{D}^{\mathsf{T}} \|^{2}$$

**Gradient:** Let  $\mathfrak{Z} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$  $\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{A}} = (\mathbf{W}_{(1)} * \mathbf{Z}_{(1)} - \mathbf{W}_{(1)} * \mathbf{X}_{(1)})(\mathbf{C} \odot \mathbf{B}) - \mathbf{Y}\mathbf{D} + \mathbf{A}\mathbf{D}^{\mathsf{T}}\mathbf{D}$  $\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{B}} = (\mathbf{W}_{(2)} * \mathbf{Z}_{(2)} - \mathbf{W}_{(2)} * \mathbf{X}_{(2)})(\mathbf{C} \odot \mathbf{A})$ 

 $\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{C}} = (\mathbf{W}_{(3)} * \mathbf{Z}_{(3)} - \mathbf{W}_{(3)} * \mathbf{X}_{(3)}) (\mathbf{B} \odot \mathbf{A})$ 

 $\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{D}} = -\mathbf{Y}^{\mathsf{T}}\mathbf{A} + \mathbf{D}\mathbf{A}^{\mathsf{T}}\mathbf{A}$ 

## **CMTF: Missing Data Estimation**

#### Construct coupled data sets and set some entries to missing in one of them

```
= [30 \ 20 \ 10 \ 15];
SZ
modes = { [1 2 3], [1 4] };
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
X{1} = full(ktensor(A(1:3)));
X{2} = tensor(A{1}*A{4}');
% add noise
for i=1:2
   N{i} = tensor(randn(size(X{i})));
   X{i} = X{i} + 0.1* N{i}/norm(N{i})*norm(X{i});
end
Xorig = X;
      = round(rand(size(X{1})));
W
X{1}(find(W==0)) = NaN;
X{1}(find(W==0)) = 0;
P = length(X);
for p=1:P
    Z.object{p} = tensor(X{p});
    obj norms(p) = norm(Z.object{p});
    Z.object{p} = Z.object{p}/obj norms(p);
end
Z.miss{1} = tensor(W);
Z.miss{2} =[];
Z.modes = modes:
Z.size = sz:
```



## **CMTF: Missing Data Estimation (cont.)**

Suppose you are given the incomplete tensor coupled with the matrix and want to jointly factorize them



## Application<sub>6</sub> (Metabolomics): Data fusion can improve missing data estimation performance

[Acar, Rasmussen, Savorani, Næs, Bro, 2013]

We have plasma samples measured using different analytical techniques, i.e., NMR and Fluorescence Spectroscopy.



### Application, (Activity Recommendation): Data fusion can improve missing data estimation performance



## Data fusion can also handle structured missing data effectively



For coupled analysis using different tensor models and loss functions, see [Ermis, Acar, Cemgil, 2015]

## Summary

### **Coupled Matrix Factorizations**

How to jointly factorize matrices using coupled matrix factorizations How to handle and estimate missing entries using CMF **Algorithms:** Gradient-based optimization algorithms

### **Coupled Matrix and Tensor Factorizations**

How to jointly factorize matrices and higher-order tensors How to jointly analyze data sets with missing entries and estimate missing data, with applications from:

### Metabolomics

### **Recommender systems**

CMTF enables us to handle and estimate structured missing data **Algorithms:** Gradient-based optimization algorithms







### **MATLAB functions**

**cmtf-opt:** Fitting CMF and CMTF to coupled data sets **(from CMTF Toolbox)** [Fac, FacInit, output] = cmtf\_opt(X, R, 'init', 'random', 'alg','ncg', 'alg\_options', options);

$$\min_{A,B,C} \left\| \mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}} \right\|^{2} + \left\| \mathbf{Y} - \mathbf{A}\mathbf{C}^{\mathsf{T}} \right\|^{2} \qquad \min_{A,B,C,D} \left\| \mathbf{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \right\|^{2} + \left\| \mathbf{Y} - \mathbf{A}\mathbf{D}^{\mathsf{T}} \right\|^{2}$$
$$\min_{A,B,C} \left\| \mathbf{W} * (\mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}}) \right\|^{2} + \left\| \mathbf{Y} - \mathbf{A}\mathbf{C}^{\mathsf{T}} \right\|^{2} \qquad \min_{A,B,C,D} \left\| \mathbf{W} * (\mathbf{X} - [\![\mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \right\|^{2} + \left\| \mathbf{Y} - \mathbf{A}\mathbf{D}^{\mathsf{T}} \right\|^{2}$$

Model: Optimization problem

**Tensor & CMTF** Toolbox Computation of the function value and the gradient

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \left\| \mathbf{X} - \mathbf{A}\mathbf{B}^{\mathsf{T}} \right\|^{2} + \frac{1}{2} \left\| \mathbf{Y} - \mathbf{A}\mathbf{C}^{\mathsf{T}} \right\|^{2} \nabla f = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{a}_{1}} \\ \vdots \\ \frac{\partial f}{\partial \mathbf{b}_{1}} \\ \vdots \\ \frac{\partial f$$

 $\left[\frac{\partial f}{\partial f}\right]$ 

**Poblano Toolbox** Gradient-based optimization algorithms

ncg, lbfgs, ...

## Limitation: CMTF assumes that all factors are shared!

In real applications, coupled data sets often have both **shared** and **unshared** factors. However, the CMTF formulation focuses on modeling only shared factors and fails in the presence of both shared/unshared factors.



## Limitation: CMTF fails to identify shared/unshared factors!

In real applications, coupled data sets often have both **shared** and **unshared** factors. However, the CMTF formulation focuses on modeling only shared factors and fails in the presence of both shared/unshared factors.



**Example 2: One shared and one unshared component in each data set**  $[\lambda_1 \ \lambda_2 \ \lambda_3] = [0 \ 1 \ 1]$ 

# Fails to identify shared and unshared components!



 $[\sigma_1 \ \sigma_2 \ \sigma_3] = [1 \ 0 \ 1]$ 

## **ACMTF (Advanced CMTF): Structure-Revealing CMTF**

#### [Acar et al., BMC Bioinformatics, 2014]

We reformulate the coupled matrix and tensor factorization problem by having factor matrices with unit norm columns and explicitly representing the weights of rank-one components in the formulation. Through modeling constraints/penalties, we let the model identify shared/unshared components.

$$\mathbf{Y} \approx \mathbf{A}\Sigma\mathbf{D}^{\mathsf{T}} \qquad \mathbf{Y} \qquad \mathbf{X} \approx [\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}]$$

$$\mathbf{Y} \approx \begin{bmatrix} \sigma_{1} & \mathbf{d}_{1} & \mathbf{d}_{R} & \mathbf{d}_{R} \\ \mathbf{u} & \mathbf{u} & \mathbf{u} \\ \mathbf{u} & \mathbf{u} \\ \mathbf{u} & \mathbf{u} \\ \mathbf{u} & \mathbf{u} \\ \mathbf{u}$$

## **ACMTF: Unconstrained Optimization**

### **Optimization Problem:**

$$\min_{A,B,C,D,\Sigma,\lambda} \| \mathcal{X} - [\lambda; A, B, C] \|^2 + \| \mathbf{Y} - A\Sigma \mathbf{D}^{\mathsf{T}} \|^2 + \beta \| \lambda \|_1 + \beta \| \sigma \|_1$$
  
s.t.  $\| \mathbf{a}_r \|_2 = \| \mathbf{b}_r \|_2 = \| \mathbf{c}_r \|_2 = \| \mathbf{d}_r \|_2 = 1$ , for  $r = 1, ..., R$ .

**Define the objective function:** 

Add as quadratic penalty terms

$$f(\lambda, \Sigma, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \| \mathbf{\mathcal{X}} - [\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}] \|^2 + \| \mathbf{Y} - \mathbf{A}\Sigma\mathbf{D}^{\mathsf{T}} \|^2 + \beta \| \lambda \|_1 + \beta \| \sigma \|_1 + \dots$$

**Smooth Approximation:** 

Replace sparsity penalties with differentiable approximations

$$f(\lambda, \Sigma, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \| \mathfrak{X} - [\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}] \|^2 + \| \mathbf{Y} - \mathbf{A}\Sigma\mathbf{D}^{\mathsf{T}} \|^2 + \beta \sum_{r=1}^{R} \sqrt{\lambda_r^2 + \epsilon} + \beta \sum_{r=1}^{R} \sqrt{\sigma_r^2 + \epsilon} + \dots$$

Compute the gradient and pick a first order optimization method

## Sparsity penalties enable us to capture the true structure



## **ACMTF: Example**

#### Construct coupled data sets with shared/unshared factors

```
sz = [30 20 10 15];
modes = {[1 2 3],[1 4]};
R = 3;
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(1,:) = [1 0 1];
lambda(2,:) = [1 1 0];
X{1} = full(ktensor(lambda(1,:)',A(1:3)));
X{2} = tensor(A{1}*diag(lambda(2,:))*A{4}');
```

```
% add noise
```

```
for i=1:2
    N{i} = tensor(randn(size(X{i})));
    X{i} = X{i} + 0.1* N{i}/norm(N{i})*norm(X{i});
end
P = length(X);
```

```
for p=1:P
Z.object{p} = tensor(X{p});
obj_norms(p) = norm(Z.object{p});
Z.object{p} = Z.object{p}/obj_norms(p);
end
Z.modes = modes;
Z.size = sz;
```



 $[\lambda_1 \ \lambda_2 \ \lambda_3] = [0 \ 1 \ 1]$ 

 $[\sigma_1 \ \sigma_2 \ \sigma_3] = [1 \ 0 \ 1]$ 





## ACMTF: Example (cont.)

#### Fit ACMTF using one of the first-order optimization algorithms

```
options = ncg('defaults');
options.MaxFuncEvals = 10000;
options.MaxIters = 10000;
options.StopTol = 1e-8;
options.RelFuncTol = 1e-8;
options.DisplayIters = 50;
% fit ACMTE-OPT
[Fac, FacInit, out] = acmtf_opt(Z,R,'alg', 'ncg','alg_options',options,'init','random', 'beta_cp',1e-3, 'beta_pca', 1e-3);
% check weights for shared/unshared factors
lrec = zeros(P,R);
for p =1:P
    temp = normalize(Fac{p});
   lrec(p,:) = temp.lambda;
end
                                                                                                             \mathfrak{X}
% lrec(1,:)*obj norms(1) will adjust the weights and correspond to
% original weights. Similarly, compute lrec(2,:)*obj norms(2) to get
% weights comparable with lambda(2,:)
                                                                                                                           ACMTF:
                                                                                                                    f(\widehat{\lambda},\widehat{\Sigma},\widehat{\mathbf{A}},\widehat{\mathbf{B}},\widehat{\mathbf{C}},\widehat{\mathbf{D}})
% compare the factors with the original factor matrices
for i=1:3
    corr(Fac{1}.U{i},A{i})
end
corr(Fac{2}.U{2},A{4})
```

## CMF models that focus on identification of shared/unshared factors

GSVD (Generalized Singular Value Decomposition) [Van Loan, 1976; Paige & Saunders, 1981] and Adapted GSVD [Van Deun et al., 2012]

Given 
$$\mathbf{X} \in \mathbb{R}^{M \times N}$$
  $(M \ge N)$ ,  $\mathbf{Y} \in \mathbb{R}^{P \times N}$   $(P \ge N)$ , let  $\mathbf{Z} = [\mathbf{X}^{\mathsf{T}} \mathbf{Y}^{\mathsf{T}}]^{\mathsf{T}}$ ,  $R = \operatorname{rank}(\mathbf{Z})$   
diagonal  
 $\mathbf{X} = \bigcup_{\mathbf{VSW}} \mathbf{W}$   
 $\mathbf{Y} = \bigcup_{\mathbf{VSW}} \mathbf{V}$   
 $\mathbf{S} = \operatorname{diag}(c_1, \dots c_N) \in \mathbb{R}^{M \times N}$   
 $\mathbf{S} = \operatorname{diag}(s_1, \dots s_N) \in \mathbb{R}^{P \times N}$   
 $c_i^2 + s_i^2 = 1$   
the relation used to

#### orthogonal

JIVE (Joint and Individual Variation Explained) [Lock et al., 2013]

Given 
$$\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, ..., \mathbf{X}^{(M)} \in \mathbb{R}^{I_m \times N}$$
, for  $m = 1, 2, ..., M$   
individual  
 $\mathbf{X}^{(1)} = \mathbf{U}^{(1)}\mathbf{S} + \mathbf{W}^{(1)}\mathbf{S}^{(1)} + \mathbf{R}^{(1)}$   
 $\mathbf{X}^{(M)} = \mathbf{U}^{(M)}\mathbf{S} + \mathbf{W}^{(M)}\mathbf{S}^{(M)} + \mathbf{R}^{(M)},$ 

the relation used to identify shared/unshared components in genome-scale expression

data sets [Alter et al., 2003]

Assumes that joint and individual parts are orthogonal!

#### BIBFA (Bayesian Interbattery Factor Analysis) [Klami et al., 2013]

Given  $X^{(m)}$ , for m=1, 2. BIBFA relies on the following probabilistic model:

## Case 1 (Easy): All methods can capture the underlying factors well!

[Acar, Bro, and Smilde, 2015]



many components similar to the shared component

## Case 1 (Easy): All methods can capture the underlying factors well!



## Case 2 (Moderate): BIBFA and ACMTF are competitive!



# Case 3 (Difficult): Increase noise level and matrix factorization-based approaches all fail!





## Uniqueness: There will be cases where ACMTF is not unique!





2

Components

3

4

1

Two unshared components in the matrix span the same subspace in different runs returning the same function value.

Since they cannot be recovered uniquely, though, they cannot be compared using the match score anymore.

# **Application<sub>8</sub> (Chemometrics): Joint analysis of measurements from multiple platforms has the potential to enhance biomarker discovery**

#### [Acar et al., BMC Bioinformatics, 2014]

**Goal:** To identify shared/unshared factors in each data set

**Data:** 29 mixtures measured using DOSY (diffusion-ordered) NMR spectroscopy and LC-MS. Mixtures are prepared using the following five chemicals:

- Val-Try-Val
- Trp Gly
- Phe
- Maltoheptaose

Can only be

captured by NMR

• Propanol

Four chemicals are expected to show up in both types of measurements!





## **ACMTF can capture the chemicals**



## ACMTF can capture the chemicals and the design

$$\begin{split} \min_{A,B,C,D,\Sigma,\lambda} \| \mathcal{X} - [\lambda; A, B, C] \|^2 + \| Y - A\Sigma D^T \|^2 + \beta \|\lambda\|_1 + \beta \|\sigma\|_1 \\ \text{s.t.} \|a_r\|_2 = \|b_r\|_2 = \|c_r\|_2 = \|d_r\|_2 = 1, \text{ for } r = 1, ..., R \\ Y \approx a_1 \quad d_1 + \dots + \sigma_6 \quad d_6 \quad \text{peaks} \quad Y_{\text{LC-MS}} \quad X_{\text{NMR}} \quad \sigma_7 \quad \lambda_r \\ \hline \mathbf{U}_{\text{LC-MS}} \quad \mathbf{U}_{\text{LC-MS}}$$

30

10

20 Mixtures
## Matrix factorization-based data fusion models fail!



Recently more studies on structurerevealing CMTF models, i.e., Khan et al., Machine Learning, 2016; Farias et al., IEEE Transactions on Signal Processing, 2016.



JIVE

20



Mixtures

GSVD



[Acar, Bro, Smilde, 2015]





BIBFA



# Different initializations may lead to different solutions!



The minimum function value:  $f(\hat{\lambda}, \hat{\Sigma}, \hat{A}, \hat{B}, \hat{C}, \hat{D}) = 0.0134$ 

Out of 250 runs with random initializations + one svd-based initialization, we get the minimum function value 24 times.



# **ACMTF** reveals chemicals visible in different platforms!

Shared

Goal: To identify shared/unshared factors in each data set

**Data:** Mixtures measured using DOSY (diffusion-ordered) NMR spectroscopy and LC-MS. Mixtures are prepared using the following five chemicals:



 $\mathfrak{X} \approx \llbracket \lambda; \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$  $\mathbf{Y} \approx \mathbf{A} \mathbf{\Sigma} \mathbf{D}^\mathsf{T}$  $\mathfrak{Z} \approx \llbracket \gamma; \mathbf{A}, \mathbf{E}, \mathbf{F} \rrbracket$ 



# Application<sub>9</sub> (Neuroscience): Joint analysis of signals from multiple neuroimaging modalities promises to provide a better understanding of the brain function

The complexity of the human brain often necessitates the use of **multiple** neuroimaging techniques to better understand neural activities. Neuroimaging techniques provide **complementary** information at different scales. **Functional Magnetic** 

**Resonance Imaging (fMRI)** 

EEG (electroencephalography) T5.01 F4-C4 C4-P4 P4-T2 Structural Fp2-F F8-T4 T4-T6 T6-T2 MRI (sMRI) electrodes voxels voxels **fMRI** subjects EEG sMRI

time samples

Can we jointly analyze these data sets and capture functional/structural patterns that differ between healthy controls and patients suffering from a disorder?

# Joint Independent Component Analysis (jICA) has been a popular approach to multi-modal neuroimaging data fusion

Single Electrode:



**Multiple Electrodes:** 

time

 $\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \dots & \mathbf{X}_K & \mathbf{Y} \end{bmatrix} = \mathbf{A}\mathbf{S}$ 

[Swinnen et al., EUSIPCO, 2014]



However, these approaches do not take into account the potential multi-linear structure of the EEG data.



### Multi-way structure can be preserved using Coupled Matrix and Tensor Factorizations



#### **Coupled tensor factorizations in neuroscience:**

EEG & magnetoencephalography (MEG) [Becker et al., EUSIPCO, 2012; Naskovska et al., EUSIPCO, 2017] EEG & Gaze data [Rivet et al., IEEE EMBC, 2015] EEG & fMRI:

Single-subject data coupled in spatial mode [Karahan et al., *Proceedings of the IEEE*, 2015] Data from multiple subjects coupled in the subjects mode modeled using CMTF [Hunyadi et al., *EUSIPCO*, 2016]

# Our Approach to Multi-modal Neuroimaging Data Fusion Calhoun, Adali, 2017]



# **Construction of Coupled Data sets: Feature Extraction**



Standard

1 kHz

Target





Standard

0.5 kHz

Standard



fMRI time series images

Task-related contrast images







sMRI images

Segmented gray matter images



Event related potential (ERP)

# **Experimental Setting**

#### Models:

### **CP Model to analyze multi-channel EEG**

**3 electrodes :** Cz, Pz, Fz **11 electrodes :** AF3, AF4, Fz, T7, C3, Cz, C4, T8, Pz, PO3, PO4 **62 electrodes** excluding VEOG/HEOG

**Structure-revealing CMTF** to jointly analyze EEG-fMRI

Structure-revealing CMTF vs. Joint ICA to jointly analyze EEG-fMRI-sMRI

Algorithms: Multiple random starts are used to deal with the local minima problem CP: CP-OPT (Tensor Toolbox) Structure-revealing CMTF: ACMTF-OPT (CMTF Toolbox) Joint ICA: ICA (Entropy Bound Minimization [Li & Adali, IEEE Trans Signal Process, 2010])

#### We study

- which electrodes to choose
- added value of joint analysis
- structure-revealing CMTF vs. jICA





# **EEG only: CP captures statistically significant factors**



Number of components: R = 3Significant components identified using attwo-sample *t*-test on the columns of A



Meaningful components, in particular, Comp1: N2-P3 transition

Highest statistical significance with 11 electrodes





2.4 x 10<sup>-4</sup>

# EEG coupled with fMRI: Structure-revealing CMTF captures meaningful temporal & spatial patterns with high resolution



Number of components: R = 10

Significant components identified using a two-sample *t*-test on the columns of A





#### Meaningful components

**Comp5:** *P2/P3 peaks in EEG, fMRI activation in the superior parietal and visual cortex* **Comp7:** *N2 -P3 transition in EEG, DMN (Default Mode Network) activation in fMRI* 

# EEG coupled with fMRI & sMRI: Structure-revealing CMTF captures meaningful temporal & spatial patterns + structure information!



# **EEG coupled with fMRI & sMRI: Joint ICA also captures components differentiating between patients and controls!**



- Subject coefficients are the same in all modalities!
- Rows of S are assumed to be statistically independent

Number of components: R = 10

**Comp1:** N2 peak in EEG, fMRI activation in the superior parietal cortex, gray matter concentration in the cerebellum

**Comp2:** *N1, P2, N2 peaks and N2-P3 transition in EEG, fMRI activation in the motor cortex and temporal lobe, gray matter concentration in the cerebellum* 



### Clustering measures demonstrate the added value of fMRI and clearly the performance difference between ACMTF and joint ICA

$$\begin{split} \boldsymbol{\mathfrak{X}} \approx \llbracket \boldsymbol{\lambda}; \mathbf{A}_{cp}, \mathbf{B}, \mathbf{C} \rrbracket & \boldsymbol{\mathfrak{X}} \approx \llbracket \boldsymbol{\lambda}; \mathbf{A}_{acmtf}, \mathbf{B}, \mathbf{C} \rrbracket & [\mathbf{X}_{(1)} \, \mathbf{Y} \, \mathbf{Z}] & \approx \mathbf{A}_{jica} \mathbf{S} \\ & \mathbf{Y} \approx \mathbf{A}_{acmtf} \boldsymbol{\Sigma} \mathbf{D}^{\mathsf{T}} \\ & \mathbf{Z} \approx \mathbf{A}_{acmtf} \boldsymbol{\Gamma} \mathbf{E}^{\mathsf{T}} \end{split}$$

Data Sets	Model	# of electrodes	Performance of k-means *				
Added va	lue of fMR	21	Accuracy (%)	F- score	Purity	Normalized Mutual Information	
EEG	СР	3	79	0.78	0.79	0.28	
(38 subjects)	СР	11	79	0.80	0.79	0.42	
	СР	62	79	0.80	0.79	0.42	
EEG -fMRI	ACMTF	3	82	0.76	0.82	0.31	
(38 subjects)	ACMTF	11	87	0.84	0.87	0.43	
	ACMTF	62	87	0.86	0.87	0.47	
EEG- fMRI - sMRI	ACMTF	11	88	0.81	0.88	0.43	
(32 subjects)	Joint ICA	11	81	0.77	0.81	0.34	

#### \*best performance out of all possible combinations of components

#### ACMTF can be formulated as a constrained optimization problem

In order to have a flexible modeling framework, we use a general-purpose optimization solver SNOPT (Sparse Nonlinear OPTimizer) [Gill et al., *SIAM Review*, 2005]

SNOPT is designed for large constrained optimization problems with smooth nonlinear functions in the objective and constraints.

SNOPT uses a sequential quadratic programming (SQP) algorithm to minimize an augmented Lagrangian.

#### Structure-revealing CMTF model:

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D},\boldsymbol{\Sigma},\boldsymbol{\lambda}} \| \boldsymbol{\mathcal{X}} - [\![\boldsymbol{\lambda};\mathbf{A},\mathbf{B},\mathbf{C}]\!] \|^2 + \| \mathbf{Y} - \mathbf{A}\boldsymbol{\Sigma}\mathbf{D}^\mathsf{T} \|^2$$
s.t.  $\| \mathbf{a}_r \|_2 = \| \mathbf{b}_r \|_2 = \| \mathbf{c}_r \|_2 = \| \mathbf{d}_r \|_2 = 1$ 

$$\sum_{r=1}^R \lambda_r \leq \beta, \sum_{r=1}^R \sigma_r \leq \beta$$
 $\sigma_r, \lambda_r \geq 0, \text{ for } r = 1, ..., R.$ 

[Acar, Nilsson, Saunders, 2014]

SNOPT  

$$\begin{array}{l} \min_{x \in \mathbb{R}^n} \phi(x) \\ \text{s.t.} \quad l \leq \begin{pmatrix} x \\ Ax \\ c(x) \end{pmatrix} \leq u \\ \\ \text{where } c(x) \text{ indicates nonlinear} \\ \text{functions, and A is a sparse matrix.} \end{array}$$

simula

# Additional constraints are easily incorporated

In many data fusion problems, we may need the following constraints to capture the underlying structures accurately and improve the interpretation.

#### **Nonnegativity Constraints:**

$$\min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D},\boldsymbol{\Sigma},\boldsymbol{\lambda}} \| \boldsymbol{\mathcal{X}} - [\![\boldsymbol{\lambda};\mathbf{A},\mathbf{B},\mathbf{C}]\!] \|^2 + \| \mathbf{Y} - \mathbf{A}\boldsymbol{\Sigma}\mathbf{D}^\mathsf{T} \|^2$$
s.t.  $\| \mathbf{a}_r \|_2 = \| \mathbf{b}_r \|_2 = \| \mathbf{c}_r \|_2 = \| \mathbf{d}_r \|_2 = 1$ 

$$\sum_{r=1}^R \lambda_r \leq \beta, \sum_{r=1}^R \sigma_r \leq \beta$$
 $\sigma_r, \lambda_r \geq 0, b_{jr}, c_{kr}, d_{mr} \geq 0$ 

for r = 1: R, j = 1: J, k = 1: K, m = 1: M.



**Angular Constraints:** When coupled data sets are overfactored, one shared factor may be represented by two closely-correlated factors. In that case, the structure-revealing model will fail to identify shared factors accurately.

$$\begin{split} \min_{\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D},\boldsymbol{\Sigma},\boldsymbol{\lambda}} \| \, \boldsymbol{\mathcal{X}} - [\![\boldsymbol{\lambda};\mathbf{A},\mathbf{B},\mathbf{C}]\!] \, \|^2 + \| \, \mathbf{Y} - \mathbf{A}\boldsymbol{\Sigma}\mathbf{D}^\mathsf{T} \, \|^2 \\ \texttt{s.t.} \quad \| \, \mathbf{a}_r \, \|_2 = \| \, \mathbf{b}_r \, \|_2 = \| \, \mathbf{c}_r \, \|_2 = \| \, \mathbf{d}_r \, \|_2 = 1 \\ | \, \mathbf{a}_r^\mathsf{T} \mathbf{a}_p | \leq \theta, | \, \mathbf{b}_r^\mathsf{T} \mathbf{b}_p | \leq \theta, | \, \mathbf{c}_r^\mathsf{T} \mathbf{c}_p | \leq \theta, | \, \mathbf{d}_r^\mathsf{T} \mathbf{d}_p | \leq \theta \\ \sum_{r=1}^R \lambda_r \leq \beta, \sum_{r=1}^R \sigma_r \leq \beta \\ \sigma_r, \lambda_r \geq 0 \text{ for } r, p \in \{1:R\}, r \neq p. \end{split}$$



# Summary

#### Structure-revealing data fusion models

Structure-revealing CMTF (ACMTF) Reformulation of CMTF to identify shared/unshared factors Can handle missing data just like CMTF Comparison of ACMTF with matrix factorization-based structure-revealing fusion models

#### **Algorithms:**

Unconstrained optimization: Gradient-based

optimization algorithms

**Constrained optimization:** Using SNOPT, we have a more flexible modeling framework

### **Applications:**

Chemometrics/metabolomics Neuroscience

#### **MATLAB** functions

acmtf-opt: Fitting an ACMTF model to coupled data sets (from CMTF Toolbox)

[Fac, FacInit, output] = acmtf\_opt(X, R, 'init', ..., 'alg',..., 'alg\_options', ..., 'beta\_cp', ..., 'beta\_pca', ....);







# Model: Optimization problem

 $\min_{A,B,C,D,\Sigma,\lambda} \| \mathcal{X} - [\lambda; A, B, C] \|^2 + \| \mathbf{Y} - \mathbf{A}\Sigma \mathbf{D}^{\mathsf{T}} \|^2 + \beta \| \lambda \|_1 + \beta \| \sigma \|_1$ s.t.  $\| \mathbf{a}_r \|_2 = \| \mathbf{b}_r \|_2 = \| \mathbf{c}_r \|_2 = \| \mathbf{d}_r \|_2 = 1$ , for r = 1, ..., R

Tensor & CMTF Toolbox Computation of the function value and the gradient

$$f(\lambda, \Sigma, \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \| \mathbf{\mathcal{X}} - [\![\lambda; \mathbf{A}, \mathbf{B}, \mathbf{C}]\!] \|^2 + \| \mathbf{Y} - \mathbf{A} \Sigma \mathbf{D}^\mathsf{T} \|^2 + \beta \sum_{r=1}^R \sqrt{\lambda_r^2 + \epsilon} + \beta \sum_{r=1}^R \sqrt{\sigma_r^2 + \epsilon} + \alpha \sum_{r=1}^R (\| \mathbf{a}_r \| - 1)^2 + \alpha \sum_{r=1}^R (\| \mathbf{b}_r \| - 1)^2 + \alpha \sum_{r=1}^R (\| \mathbf{c}_r \| - 1)^2 + \alpha \sum_{r=1}^R (\| \mathbf{d}_r \| - 1)^2, \nabla f = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{a}_R} \\ \frac{\partial f}{\partial \mathbf{b}_R} \\ \frac{\partial f}{\partial \mathbf{c}_R} \\ \frac{\partial f}{\partial \mathbf{d}_R} \\ \frac{\partial f}{\partial \mathbf{d}_R} \end{bmatrix}$$

Poblano Toolbox Gradient-based optimization algorithms

ncg, lbfgs, ...

 $rac{\partial f}{\partial \sigma}$ 



# For what we have not discussed (rank, uniqueness, more on algorithms and constraints, other tensor models, scalability, many more applications...)

A. Smilde, R. Bro, P. Geladi. Multi-way Analysis: Applications in the Chemical Sciences, Wiley, England, 2004

- T. G. Kolda, B. W. Bader. Tensor Decompositions and Applications. SIAM Review, 51(3): 455-500, 2009
- E. Acar, B. Yener. Unsupervised Multiway Data Analysis: A Literature Survey. *IEEE Transactions on Knowledge and Data Engineering*, 21(1):6-20, 2009
- L. Grasedyck, D. Kressner, C. Tobler. A literature survey of low rank tensor approximation techniques. *GAMM-Mitteilungen*, 36: 53-78, 2013
- P. Comon. Tensors: A Brief Introduction. IEEE Signal Processing Magazine, 31(3): 44-53, 2014
- E. E. Papalexakis, C. Faloutsos, N. D. Sidiropoulos. Tensors for Data Mining and Data Fusion: Models, Applications, and Scalable Algorithms. *ACM Transactions on Intelligent Systems and Technology*, 8(2), Article 16, 2016
- N. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, C. Faloutsos. Tensor Decomposition for Signal Processing and Machine Learning. *IEEE Transactions on Signal Processing*, 65(13): 3551-3582, 2017

## Acknowledgments



CP-OPT CP-WOPT CMTF-OPT







Tamara G. Kolda, Sandia National Labs, CA

Daniel M. Dunlavy, Sandia National Labs, NM

Structure-revealing CMTF and chemometrics/metabolomics applications

Rasmus Bro, University of Copenhagen, Denmark

Age K. Smilde, University of Amsterdam, The Netherlands

Using structure-revealing CMTF for fusion of multi-modal neuroimaging data

Morten A. Rasmussen, University of Copenhagen, Denmark

Tülay Adalı, University of Maryland Baltimore County, MD



Structure-revealing CMTF using SNOPT

Michael Saunders, Stanford University, CA



Link prediction using CMTF with different loss functions and tensor models

Taylan Cemgil, Bogazici University, Turkey

Beyza Ermis, Bogazici University, Turkey

# **JODA: Joint Data Analysis for Enhanced Knowledge Discovery**

www.models.life.ku.dk/joda



#### *Literature on Coupled factorizations: Over 120 papers as of June 2018*

I. Van Mechelen, A. K. Smilde. A generic linked-mode decomposition model for data fusion. *Chemometrics an Intelligent Laboratory Systems*, 104: 83-94, 2010

E. Acar, R. Bro, A. K. Smilde. Data Fusion in Metabolomics based on Coupled Matrix and Tensor Factorizations. *Proceedings of the IEEE*, 103: 1602-1620, 2015

D. Lahat, T. Adalı, C. Jutten. Multimodal Data Fusion: An Overview of Methods, Challenges, and Prospects. *Proceedings of the IEEE*, 103: 1449-1477, 2015

# MATLAB CMTF Toolbox:



#### The MATLAB CMTF Toolbox

The MATLAB CMTF Toolbox has two different versions of the Coupled Matrix and Tensor Factorization (CMTF) approach used to jointly analyze datasets of different orders: (i) CMTF and (ii) ACMTF. First-order unconstrained optimization is used to fit both versions. The MATLAB CMTF Toolbox has the functions necessary to compute function values and gradients for CMTF and ACMTF. For the optimization routines, it uses the Poblano Toolbox. The Tensor Toolbox is also needed to run the functions in the MATLAB CMTF Toolbox.

In order to learn about the coupled models in the toolbox and see example scripts showing how to use CMTF and ACMTF, please visit the examples page.

#### What is new in Version 1.1? (Dec., 2014)

- O Compatibility with sptensor is added to make CMTF\_OPT and ACMTF\_OPT work with tensors in sptensor form.
- STESTER\_CMTF and TESTER\_ACMTF have been modified to have the option of generating data sets in dense or sparse tensor format.
- TESTER\_CMTF\_MISSING and TESTER\_ACMTF\_MISSING functions have been added to demonstrate the use of CMTF\_OPT and ACMTF\_OPT with data sets with missing entries.
- O CREATE\_COUPLED\_SPARSE function has been added to generate coupled sparse data sets using sparse factor matrices.
- For smooth approximation of the l1-terms in SCP\_FG, SCP\_WFG, SPCA\_FG, SPCA\_WFG, eps is set to 1e-8.

ſ	JabRef References output × +			l	- 0			
	www.models.life.ku.dk/~acare/DataFusion			V C Q Search	☆ 自 ♥ ♣ ♠	¢		
	van Deun, K., Smilde, A.K., van der Werl, M.J., Klers, H.A. & van Mechelen, I.	A structured overview of simultaneous component based data integration [BibTeX]	2009	BMC Bioinformatics Vol. 10, pp. 246	article	0	L	
	van Deun, K., Wilderjans, T.F., van den Berg, R.A., Antoniadis, A. & van Mechelen, I.	A flexible framework for sparse simultaneous component based data integration $[{\tt BibTeX}]$	2011	BMC Bioinformatics Vol. 12, pp. 448	article	UF	L.	
	Drumond, L.R., Diaz-Aviles, E., Schmidt-Thieme, L. & Nejdl, W.	Optimizing Multi-Relational Factorization Models for Multiple Target Relations [BibTeX]	CIKM14: Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, pp. 191-200	e inproceedings	UF	L		
le	Ermis, B., Acar, E. & Cemgil, A.T.	Link prediction in heterogeneous data via generalized coupled tensor factorization [BibTeX]	2015	Data Mining and Knowledge Discovery Vol. 29, pp. 203-236	article	UF	URL	
	Ernis, B., Acar, E. & Comgli, A.T.	Link Prediction via Generalized Coupled Tensor Factorisation [BibTeX]	2012	ECML/PKDD Workshop on Collective Learning and Inference on Structured Data	inproceedings	UF	URL URL	
	Farias, R.C., Cohen, J.E., Jutten, C. & Comon, P.	Joint decompositions with flexible couplings [BibTeX]	2015		techreport	UF		
	Groves, A.R., Beckmann, C.F., Smith, S.M. & Woolrich, M.W.	Linked independent component analysis for multimodal data fusion [BibTeX]	2011	NeuroImage Vol. 54, pp. 2198-2217	article	DC	DOI	
	Opera. S., K. Phong, D., Adam, E., Tran, T. &         Nonregative Shared Subsyste Learning and Its Application to Social Media Retrieval (IBB-RA)         20           Harshman, R.J. & Lundy, M.S.         Exba proprocessing and the extended PADAFAC model         110           Harshman, R.J. & Lundy, M.S.         PADAFAC Parallel Factor Analysis         119           Harshman, R.J. & Lundy, M.S.         PADAFAC Parallel Factor Analysis         119		2010	KDD10: Proc. 16th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, pp. 1169-1178	inproceedings		URL	
			1984	Research Methods for Multi-mode Data Analysis, pp. 216-284	incollection			
			1994	Computational Statistics and Data Analysis Vol. 18, pp. 39-72	article	UF	JRL	
	Hotelling, H.	Relations Between Two Sets of Variates	1936	Biometrika Vol. 28, pp. 321-377	article	UF	L.	

### Evrim Acar evrim@simula.no

www.models.life.ku.dk/~acare