

Environmental Multiway Data Mining

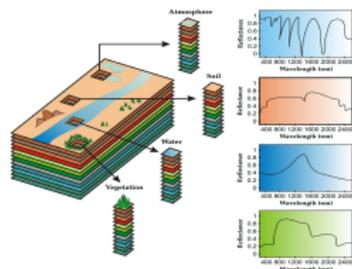


Jérémy E. Cohen, Nicolas Gillis

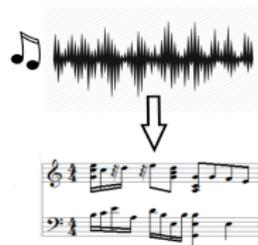
UMONS, FNRS

March 21, 2018

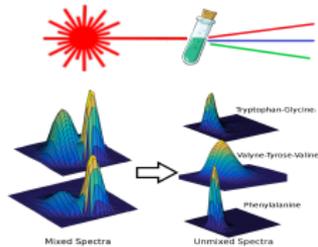
Tensors in Signal Processing



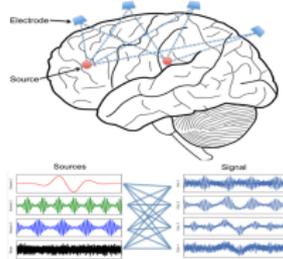
TELEDETECTION



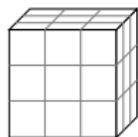
AUDIO



CHIMOMETRIE



NEUROSCIENCES

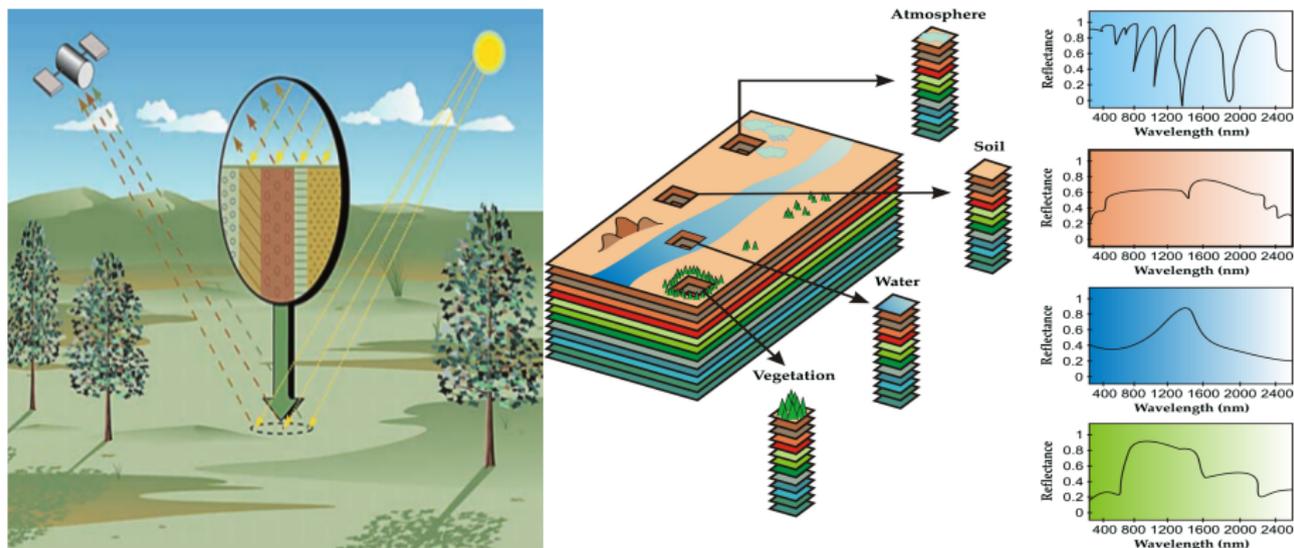


Data tensor

Ex : $\text{freq} \times \text{time} \times \text{sensor}$

Tensor model = Accounting for structure

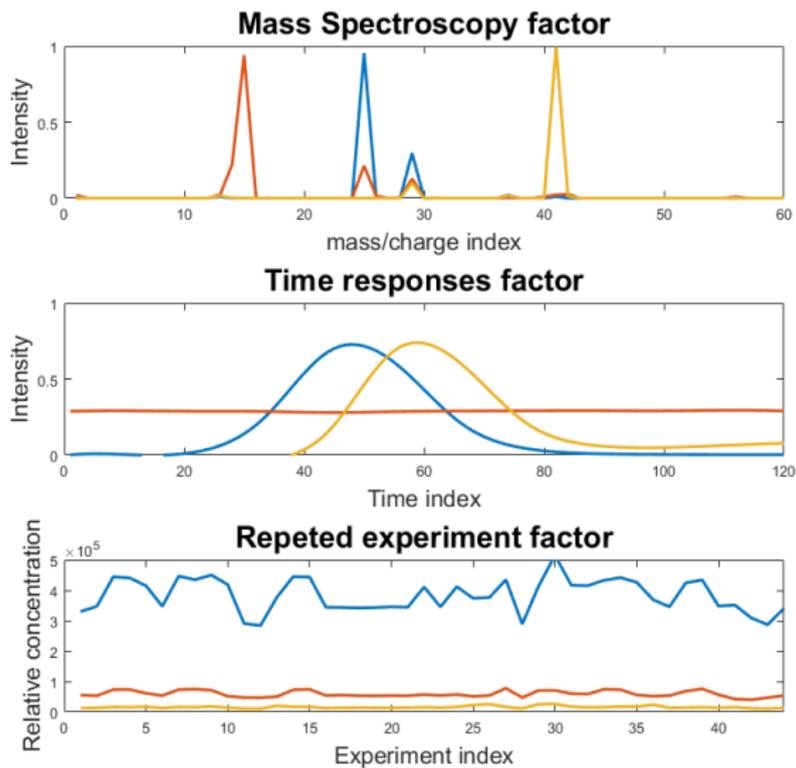
Hyperspectral imaging principle



- Each image is a **linear mixture** of various spectral signatures.
- Each material has a unique spectral response.

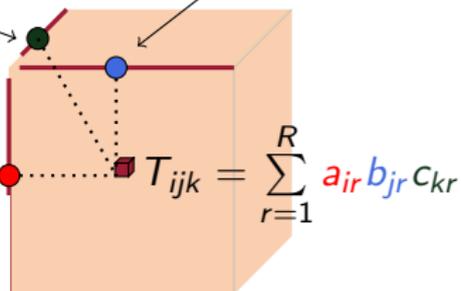
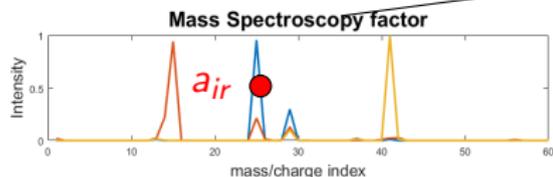
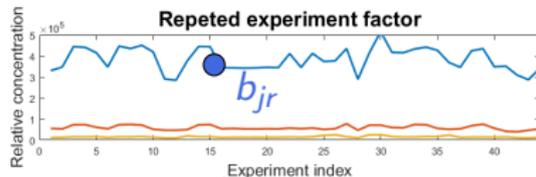
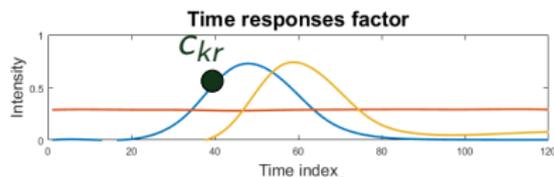
Credits for illustrations: Veganzones (left) and Bioucas (right)

Liquid Chromatography — Mass Spectroscopy component



Mass over charge intensities form a tensor

$$T(\lambda, t, k) = \sum_{r=1}^R a_r(\lambda) b_r(t) c_r(k)$$



A tool for LRA: Canonical Polyadic Decomposition

Canonical Polyadic Decomposition [Hitchcock,1927] aims at extracting all R components.

Tensor = first component + ... + R th component

- Unmixing in theory does not require additional knowledge for order 3 and more.
- For matrices, not unique if $R > 1 \rightarrow$ SVD (orthogonality), NMF (non-negativity).

CPD

$$\mathcal{T} = \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c}_1 + \dots + \mathbf{a}_R \otimes \mathbf{b}_R \otimes \mathbf{c}_R$$

$$\mathcal{T} = (\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C}) \mathcal{I}_R$$

\mathcal{T} has sizes $K \times L \times M$

\otimes is the tensor product

R is the rank of \mathcal{T} , *i.e.* smallest number of rank-one tensors spanning \mathcal{T} .

$\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R]$ has sizes $K \times R$

\bullet_i is the contraction on mode i

Tensor decomposition as an approximation problem

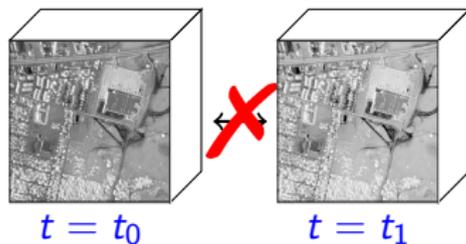
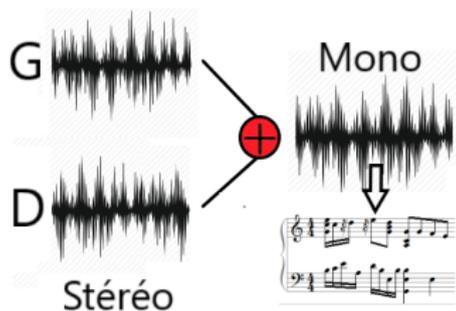
$$\begin{array}{ll} \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} & \|\mathcal{T} - (\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C})\mathcal{I}_R\| \\ \text{sub. to} & \mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathcal{C}_{A, B, C} \end{array}$$

- Non-convex in the general case but convex with respect to each block $\mathbf{A}, \mathbf{B}, \mathbf{C}$.
- Example: Non-negative Matrix Factorization with Frobenius norm

$$\begin{array}{ll} \min_{\mathbf{A}, \mathbf{B}} & \|\mathbf{M} - \mathbf{A}\mathbf{B}^T\|_F^2 \\ \text{sub. to} & \mathbf{A} \geq 0 \quad \mathbf{B} \geq 0 \end{array}$$

Challenges in tensor signal processing

Multidimensional structure **not exploited** !



Main issues:

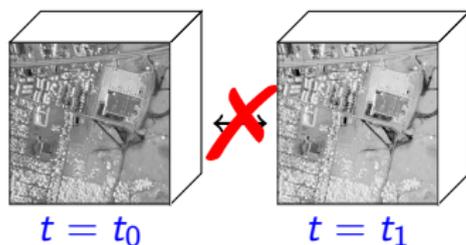
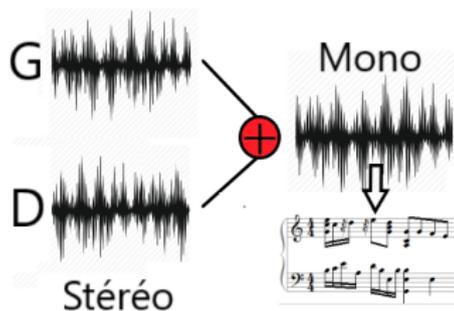
Interpretability

Tensor formalism

Constraints / Size

Challenges in tensor signal processing

Multidimensional structure **not exploited** !



Main issues:

Interpretability

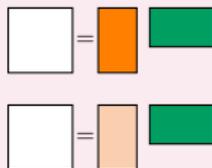
Tensor formalism

Constraints / Size

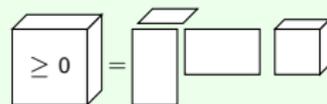
I. Informed decom.



II. Data fusion

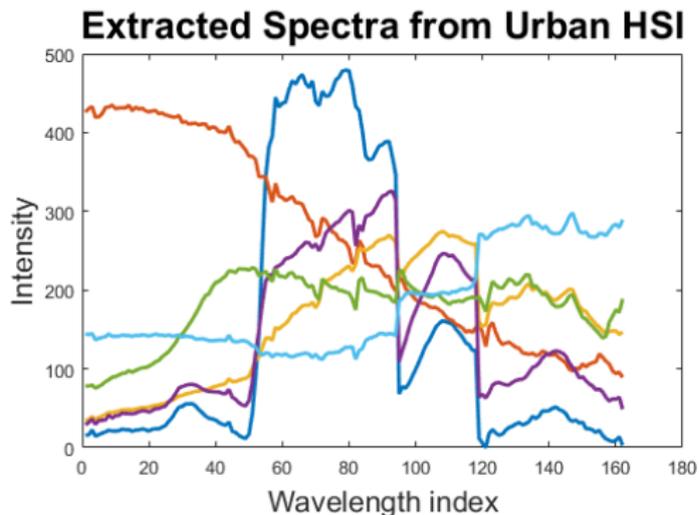


III. Optimisation

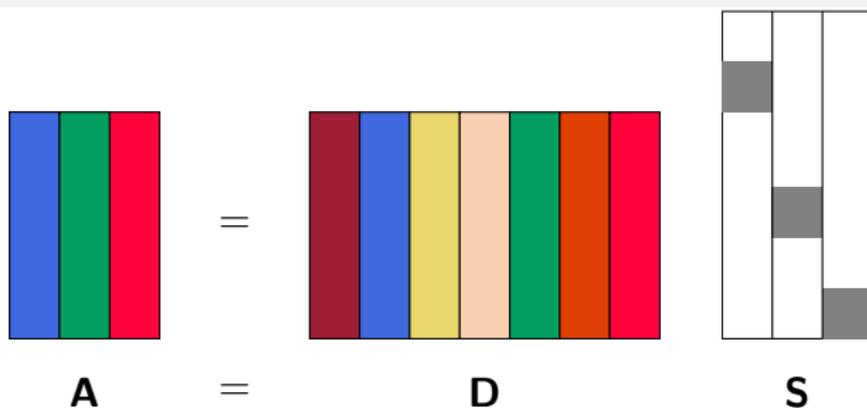


Challenge I: Identification.

Identification may be an issue



Let's choose **A** from a dictionary



$$\mathbf{X} = \mathbf{D}\mathbf{S}\mathbf{B}^T = \mathbf{D}(:, \mathcal{K})\mathbf{B}^T \text{ or } \mathcal{T} = (\mathbf{D}\mathbf{S} \otimes \mathbf{B} \otimes \mathbf{C})\mathcal{I}_R \text{ where } \|\mathbf{S}\|_{col,0} = 1$$

1. Theorem: If $\text{spark}(\mathbf{D}) > R$ and \mathcal{K} has no repetition,

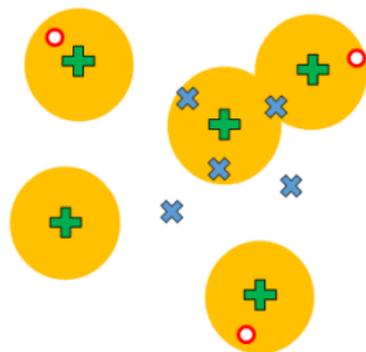
(i) if there exist $(\mathcal{K}, \mathbf{B})$ so that $\mathbf{M} = \mathbf{D}(:, \mathcal{K})\mathbf{B}$, then it is unique.

(ii) $\underset{\mathcal{K}, \mathbf{B}, \mathbf{C}}{\text{argmin}} \|\mathcal{T} - (\mathbf{D}(:, \mathcal{K}) \otimes \mathbf{B} \otimes \mathbf{C})\mathcal{I}_R\|_F^2$ existe.



Flexibility and Separability

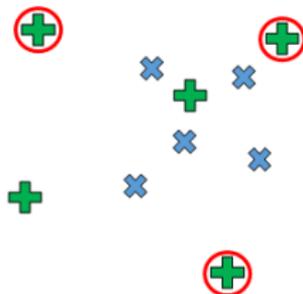
Flexibility



- × Data points
- + Atoms in D
- Factor matrix A
- Search space

$$A \approx DS$$

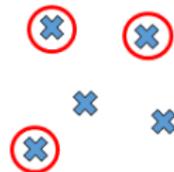
Standard case



- × Data points
- + Atoms in D
- Factor matrix A

$$A = DS$$

Separability



- × Data points
- + Atoms in D
- Factor matrix A

$$A = XS$$

Matching Pursuit ALS (works for high-order tensors)

An alternating nonnegative least squares method where \mathbf{A} , \mathcal{K} and \mathbf{B} are estimated alternatively.

Input: Initial \mathbf{A} , \mathbf{B} , \mathcal{K} (Using e.g. SPA, VCA...).

Run a few iterations of NMF.

while stopping criterion is not met,

- $\hat{\mathbf{A}} = \underset{\mathbf{A} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{A}\mathbf{B}^T\|_F^2$
- $\hat{\mathcal{K}}(i) = \underset{j}{\operatorname{argmax}} d_j^T \hat{\mathbf{A}}_i \quad \forall i \in [R]$
- $\hat{\mathbf{B}} = \underset{\mathbf{B} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{D}(:, \hat{\mathcal{K}})\mathbf{B}^T\|_F^2$

Output: Selected atoms set \mathcal{K} and abundances \mathbf{B} .

Pros and Cons

Pros

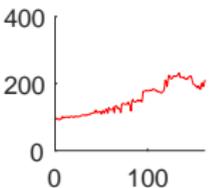
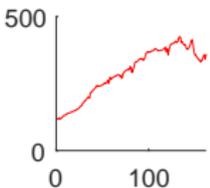
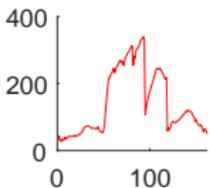
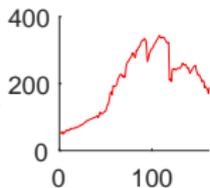
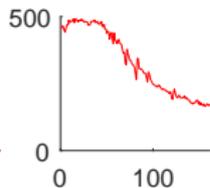
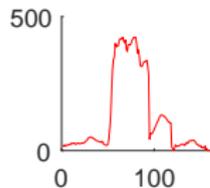
- ✓ Can be adapted to N-way arrays.
- ✓ Can be adapted for more complex estimation schemes of A and B .
- ✓ One iteration has the same complexity as geometric methods.
- ✓ Low memory requirements.
- ✓ Tries to minimize an explicit cost function.

Cons

- ✗ Very sensitive to initialization.
- ✗ No convergence proof.
- ✗ Requires the knowledge of R .

Application to Spectral Unmixing with Pure pixels

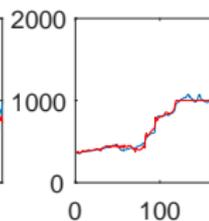
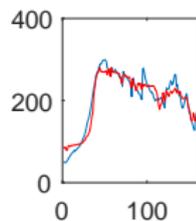
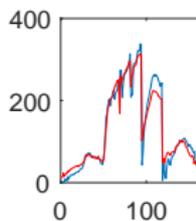
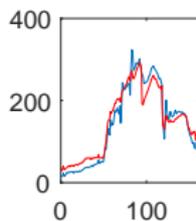
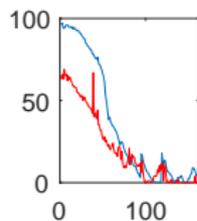
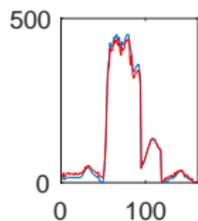
Spectra extracted exactly from the data (in red)



Spectral band index

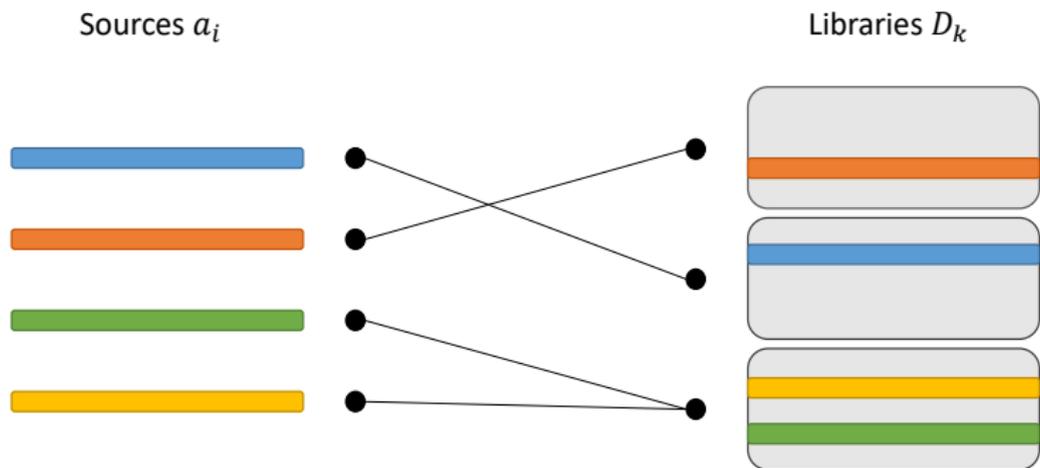
Application to Spectral Unmixing with Pure pixels

Spectra (in blue) close the data (in red)



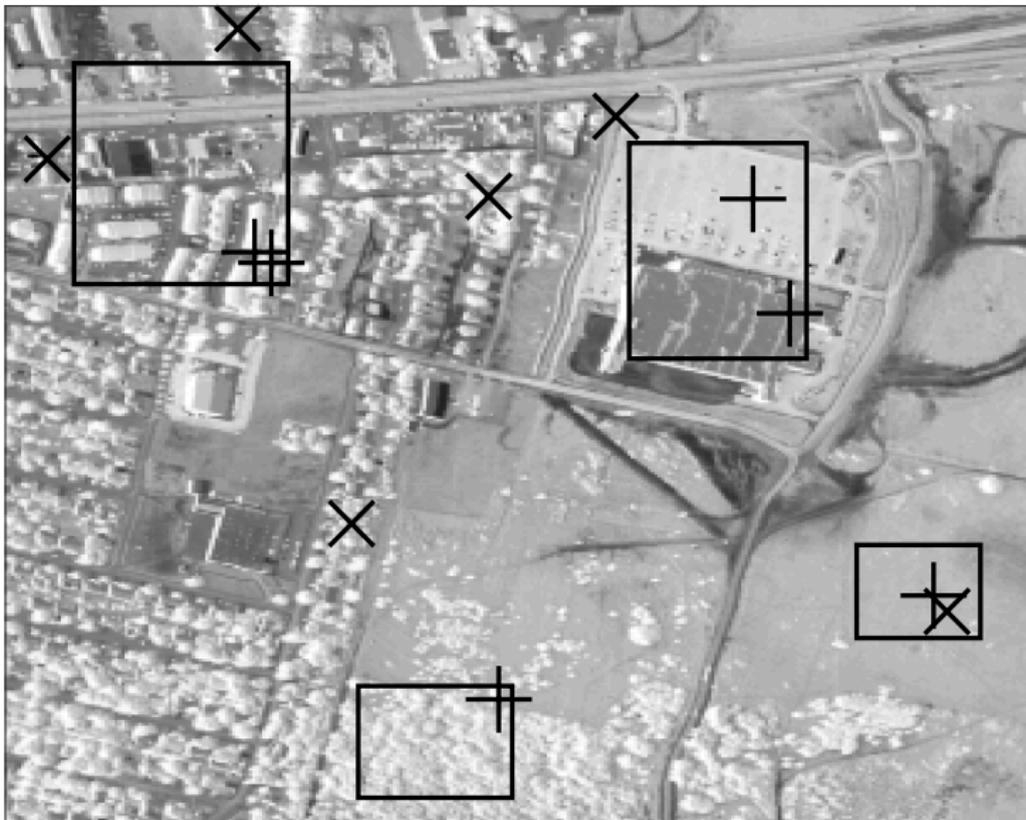
Toolbox available on my personal webpage jeremy-e-cohen.jimdo.com
[Cohen Gillis, 2017]

Multiple Dictionary for Hand-Picking Pure Pixels

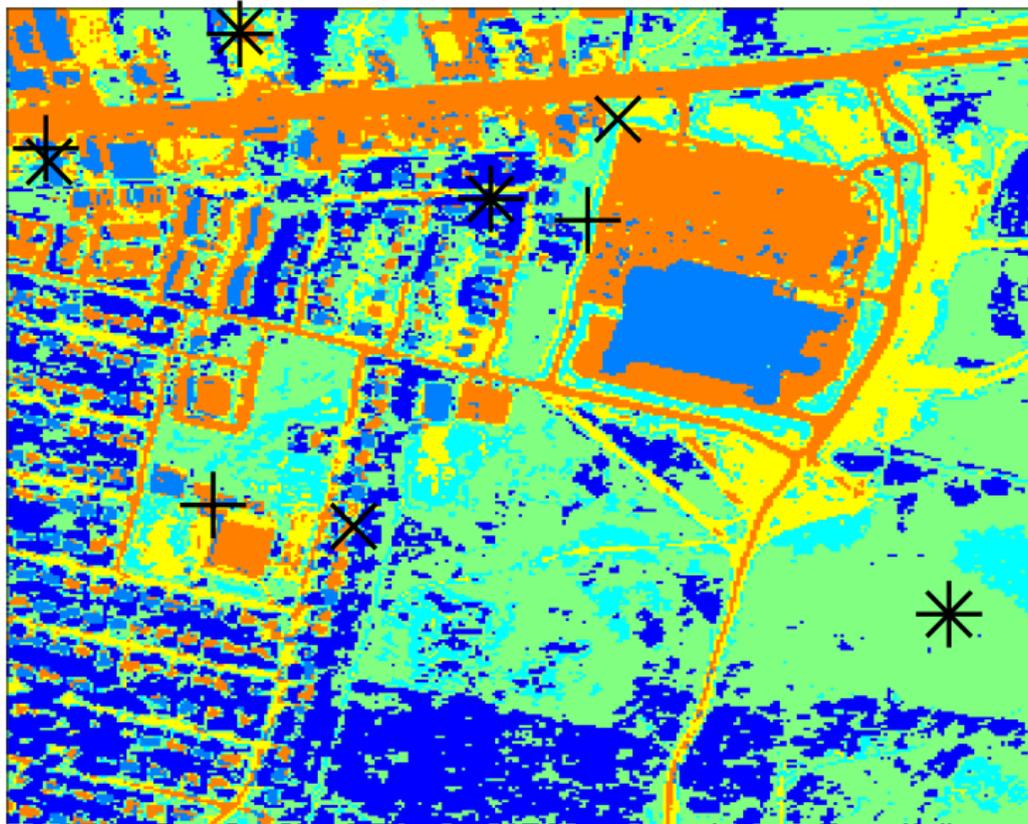


$$\mathbf{A} = [\mathbf{D}_1(:, \mathcal{K}_1), \dots, \mathbf{D}_p(:, \mathcal{K}_p)] \mathbf{\Pi}$$

Example: Supervised Multiple Dictionary learning



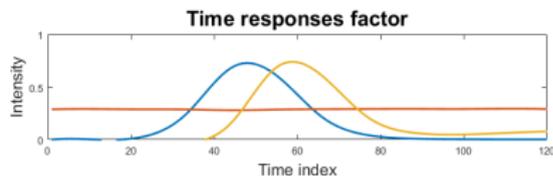
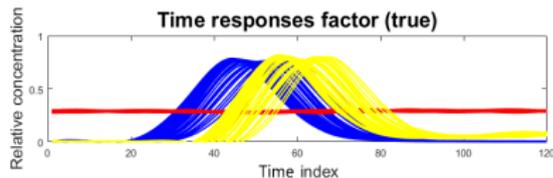
Example: Unsupervised version using segmentation



Challenge II: Subject Variability and Multimodality.

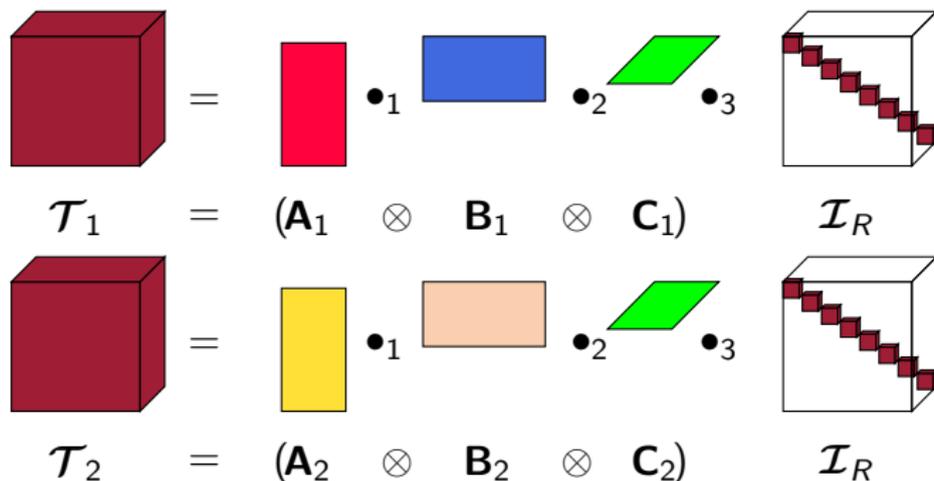
Subject Variability

$$\begin{array}{c}
 \text{[Red Square]} \\
 \mathbf{M}_1 \\
 \\
 \vdots \\
 \\
 \text{[Red Square]} \\
 \mathbf{M}_i
 \end{array}
 =
 \begin{array}{c}
 \text{[Red Rectangle]} \\
 \mathbf{A} \\
 \\
 \vdots \\
 \\
 \text{[Red Rectangle]} \\
 \mathbf{A}
 \end{array}
 \begin{array}{c}
 \text{[Diagonal Matrix]} \\
 \boldsymbol{\Sigma}_1 \\
 \\
 \vdots \\
 \\
 \text{[Diagonal Matrix]} \\
 \boldsymbol{\Sigma}_i
 \end{array}
 \begin{array}{c}
 \text{[Blue Rectangle]} \\
 \mathbf{B}_1^T \\
 \\
 \vdots \\
 \\
 \text{[Light Blue Rectangle]} \\
 \mathbf{B}_i^T
 \end{array}$$



Example: LC-MS data.

Data Fusion with tensors



Example: Fluorescence and NMR data. Often $\mathbf{C}_1 := \mathbf{C}_2$. But:

- the sampling rates can be different?
- the relation may not be trivial? Can it be learned?
- how does coupling affect the cost function?

General Framework using a Bayesian approach [Cabral Farias, Cohen,2015]

- Parameters $\theta_i = [\text{vec}(\mathbf{A}_i); \text{vec}(\mathbf{B}_i); \text{vec}(\mathbf{C}_i)]$ are random
- Known prior distribution $p(\theta_1, \dots, \theta_N)$ and likelihoods $p(\mathcal{Y}_i | \theta_i)$
- Deterministic point of vue: $\theta_i = \phi_i(\theta^*)$ for some fixed function ϕ_i .

MAP estimation under conditionnal independance

$$\arg \max_{\theta_1, \dots, \theta_N} p(\theta_1, \dots, \theta_N | \mathcal{Y}_1, \dots, \mathcal{Y}_N) = \arg \min_{\theta_1, \dots, \theta_N} \Upsilon(\theta_1, \dots, \theta_N)$$

$$\begin{aligned} \Upsilon(\theta_1, \dots, \theta_N) &= - \sum_{i=1}^N \log p(\mathcal{Y}_i | \theta_i) - \log p(\theta_1, \dots, \theta_N) \\ &= \text{data fitting terms} + \text{coupling} \end{aligned}$$

Some flexible coupled LRA models

Noisy exact coupling on \mathbf{C}_i

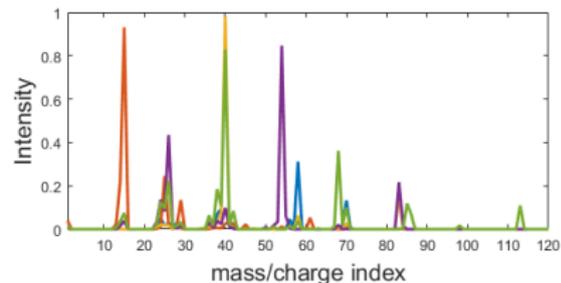
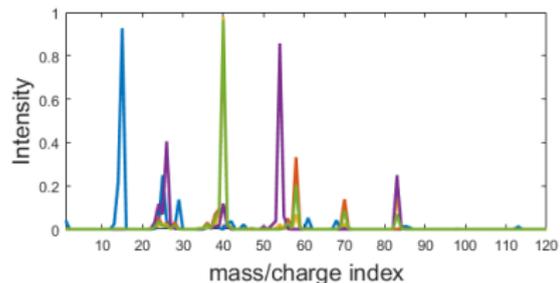
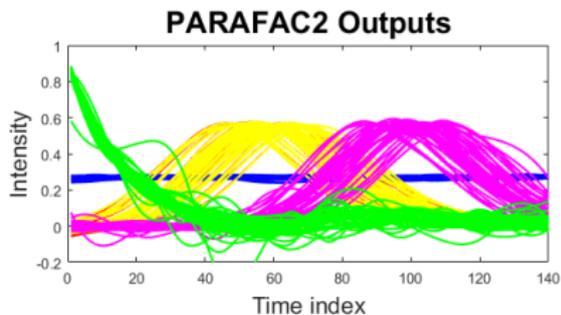
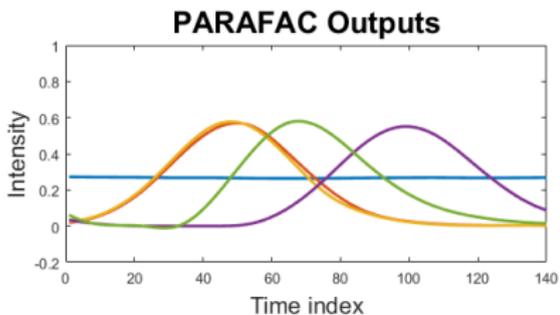
$$\forall i \in [1, M], \quad \begin{cases} \mathcal{T}_i &= (\mathbf{A}_i \otimes \mathbf{B}_i \otimes \mathbf{C}_i) \mathcal{I}_R + \mathcal{E}_i \\ \mathbf{C}_i &= \mathbf{C}^* + \mathbf{\Gamma}_i \\ \mathbf{\Gamma}_i &\sim \mathcal{N}\left(\mathbf{0}, \frac{1}{\sigma_{c,i}^2} \mathbf{I} \otimes \mathbf{I}\right) \end{cases}$$

$$\Upsilon(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, \mathbf{C}^*) = - \sum_{i=1}^N \frac{1}{\sigma_1^2} \|\mathcal{T}_i - (\mathbf{A}_i \otimes \mathbf{B}_i \otimes \mathbf{C}_i) \mathcal{I}_R\|_F^2 - \sum_{i=1}^N \frac{1}{\sigma_{ci}^2} \|\mathbf{C}_i - \mathbf{C}^*\|_F^2$$

PARAFAC2 [Harshman,1972][Bro,1999][Cohen,2018]

$$\forall i \in [1, M], \quad \begin{cases} \mathbf{M}_i &= \mathbf{A}_i \boldsymbol{\Sigma}_i \mathbf{B}_i^T + \mathbf{E}_i \\ \mathbf{A}_i &= \mathbf{A}^* \\ \mathbf{B}_i &= \mathbf{P}_i \mathbf{B}^* \\ \mathbf{P}_i^T \mathbf{P}_i &= \mathbf{I} \end{cases}$$

PARAFAC2 vs PARAFAC on LC-MS data



Many other solutions can be thought of to tackle subject variability!

Perspectives

Variability along time / sensors

- Characterize variations along time/sensors within a multiway model?
Or in a statistical manner, *i.e.* with priors on the evolution of coupled parameters?
- Applications: automatic stereo transcriptions, temporal spectral unmixing and super resolution. . .

Interactions between machine learning and multimodality

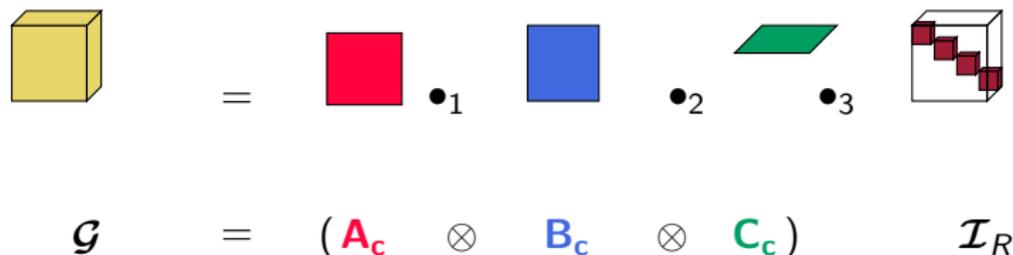
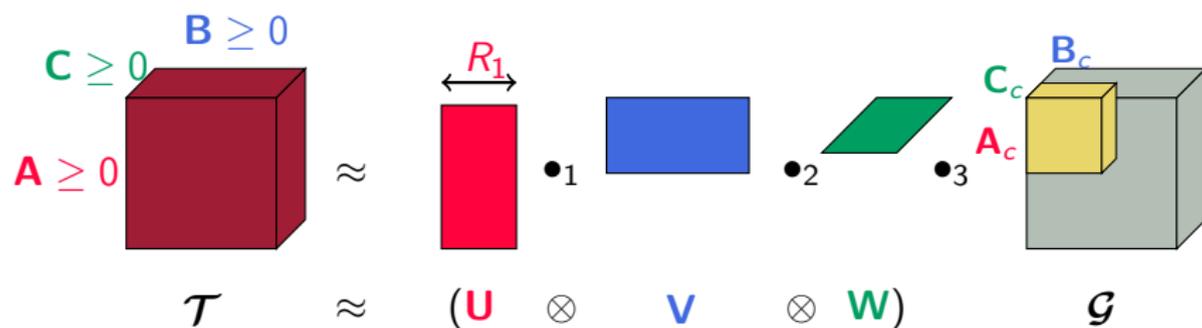
- Learning the coupling relationship
- Tensor dictionary learning

Challenge III: Constrained large decompositions.

DEGA BUREN:
ASTÉRIX LE GAULOIS
LA SERPE D'OR
ASTÉRIX ET LES GOTS
ASTÉRIX GLADIATEUR
LE TRAI DE GAULE
ASTÉRIX ET CLEOPÂTRE
LE GOMBYAT DES CHIENS
ASTÉRIX CHEZ LES NÉTOS
ASTÉRIX ET LES MORGANDS
ASTÉRIX LEGIONNAIRE
LE BOUCLIER ARMÉNIEN
ASTÉRIX AU JEU GLOUPESS
ASTÉRIX ET LE CHAUDRON
ASTÉRIX EN HESPIE
LA ZÉLANIE
ASTÉRIX CHEZ LES HÉLVÉTIENS
LE DORGANE DES BOULES
LES LAURENDS DE CÉSAR
LE DIEU NIN
ASTÉRIX ET CÉSAR
LE CASQUE DE CÉSAR
LA GRANDE TRAVERSÉE
CÉSAR ET COMPAGNIE



Unconstrained compression ...



$$\mathcal{T} \approx (\mathbf{U} \otimes \mathbf{V} \otimes \mathbf{W}) \mathcal{G} = (\mathbf{U}\mathbf{A}_c \otimes \mathbf{V}\mathbf{B}_c \otimes \mathbf{W}\mathbf{C}_c) \mathcal{I}_R$$

... but constrained decomposition!

Compressed domain NN CP:

$$\begin{array}{ll} \min_{\mathbf{A}_c, \mathbf{B}_c, \mathbf{C}_c} & \|\mathcal{G} - (\mathbf{A}_c \otimes \mathbf{B}_c \otimes \mathbf{C}_c)\mathcal{I}\| \\ \text{sub. to} & \widehat{\mathbf{U}}\mathbf{A}_c, \widehat{\mathbf{V}}\mathbf{B}_c, \widehat{\mathbf{W}}\mathbf{C}_c \geq 0 \end{array}$$

Issue

Solution

Easy unconstrained/difficult constrained Unconstrained solution \rightarrow projectionDifficult exact projection $\widehat{\mathbf{U}}\mathbf{A}_c$

Approximate projection

Approximate projection and PROCO-ALS

Approximate projection Π :

Given Least Squares update $\hat{\mathbf{A}}_c$

① Decompression: $\hat{\mathbf{A}} := \hat{\mathbf{U}}\hat{\mathbf{A}}_c$

② Projection: $\hat{\mathbf{A}} := [\hat{\mathbf{A}}]^+$

③ Compression: $\hat{\mathbf{A}}_c := \hat{\mathbf{U}}^T\hat{\mathbf{A}}$

$$\Pi [\hat{\mathbf{A}}] = \mathbf{U}^T[\mathbf{U}\mathbf{A}_c]^+$$

Projected and compressed framework (PROCO) [Cohen,2014]

Other possible algorithms and related problems

- PROCO-ALS [Cohen,2014], Compressed-AOADMM [Cohen,2016]

$$\begin{aligned} & \text{minimize } \|\widehat{\mathcal{G}} - (\mathbf{A}_c \otimes \mathbf{B}_c \otimes \mathbf{C}_c) \mathcal{I}_R\|_F^2 \\ & \text{w.r.t. } \mathbf{A}_c, \mathbf{B}_c, \mathbf{C}_c \\ & \text{s.t. } \widehat{\mathbf{U}} \mathbf{A}_c \succeq 0 \end{aligned}$$

- Tensorlab 3.0 [Vervliet,2016]

$$\begin{aligned} & \text{minimize } \|\left(\widehat{\mathbf{U}} \otimes \widehat{\mathbf{V}} \otimes \widehat{\mathbf{W}}\right) \widehat{\mathcal{G}} - (\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C}) \mathcal{I}_R\|_F^2 \\ & \text{w.r.t. } \mathbf{A}, \mathbf{B}, \mathbf{C} \\ & \text{s.t. } \mathbf{A} \succeq 0 \end{aligned}$$

- AOADMM [Huang,2015], FastNNLS [Bro,1997], ANLS

$$\begin{aligned} & \text{minimize } \|\mathcal{T} - (\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C}) \mathcal{I}_R\|_F^2 \\ & \text{w.r.t. } \mathbf{A}, \mathbf{B}, \mathbf{C} \\ & \text{s.t. } \mathbf{A} \succeq 0 \end{aligned}$$

Application in Fluorescence Spectroscopy

Fluorescence spectroscopy data: excitation spectra
emission spectra
mixtures

multimodal chemometrics data set from Acar *et al*¹

Description

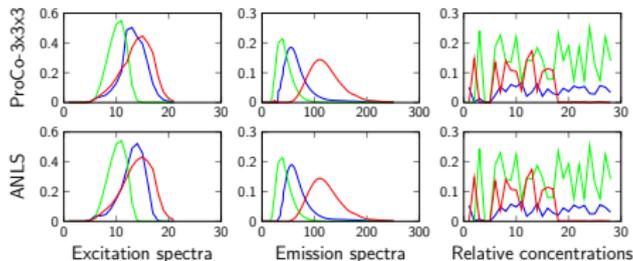
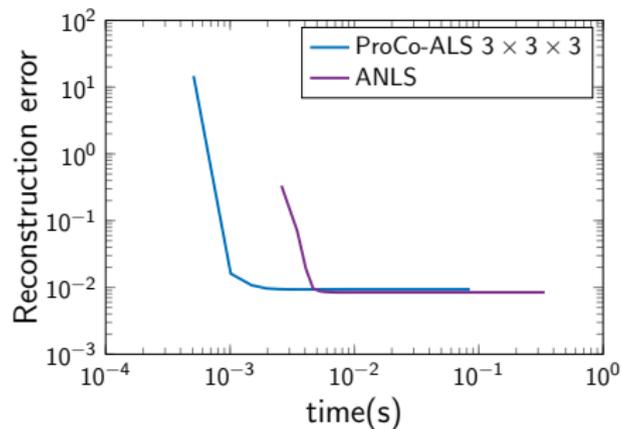
5 compounds: Valine-Tyrosine-Valine (Val), Tryptophan- Glycine (Gly), Phenylalanine (Phe), Maltoheptaose (Mal) and Propanol (Pro)

Nb. of excitation wave lengths	21 (A)
Nb. of emission wave lengths	251 (B)
Nb. of Mixtures	28 (C)
Missing values	30% (replaced by zeros)

¹E. Acar, A.J. Lawaetz, M.A. Rasmussen, and R. Bro. Structure-revealing data fusion model with applications in metabolomics. In Conf. Proc. IEEE Eng. Med. Biol. Soc., pages 6023– 6026. IEEE, 2013

Application to Fluorescence Spectroscopy

ANLS (nonnegative) and ProCo-ALS



Conclusions and Perspectives

Studied topics

- Identification through known dictionaries.
- Multimodality and Subject Variability in matrix/tensor low rank factorization models for chemometrics/neuroimaging.
- Constrained tensor compression and decomposition, especially in the context of nonnegativity.

Things I am (or would like to be) working on

- Flexible dictionary constraints, tensor dictionary learning.
- Data fusion for temporal series of hyperspectral images.
- Multispectral/Hyperspectral fusion for spectral unmixing.
- Audio source separation with tensor models, which calls for new tensor decomposition models and non-euclidean error metrics.

Thank you for your attention!



State-of-the-art (non-exhaustive)

- **Continuous approaches**

Lasso, GLUP [Ammanouil 2014], FGNSR [Gillis 2016]

+ Robust, optimization criterion. – Slow.

- **Greedy/Non-iterative method**

- Geometric algorithms (pure pixel hypothesis)

N-FINDR [Winter 1999], VCA [Nascimento 2005], SPA [Gillis 2014, Businger Golub 1965]

- Matching pursuit approaches

SDSOMP[X.Fu 2013, Tropp 2006]

+ Fast – Not robust, No explicit criterion

- **Pixel-wise brute force algorithms**

MESMA [Roberts 1998], MESLUM, AUTOMCU, AMUSES [Degerickx 2017]

+ Flexible – No low rank property, Slow.

- **Statistical methods**

Experiment on the URBAN HSI

	$r = 6$		$r = 8$	
	Time (s.)	Rel. err.	Time (s.)	Rel. err.
RAND-wo	0.00	7.87	0.00	11.66
d-RAND-wo	22.46 (13)	5.09	34.87 (18)	5.35
RAND-av	0.02	11.51	0.02	9.60
d-RAND-av	23.91 (13)	4.65	30.77 (15)	4.65
RAND-be	0.00	13.77	0.00	5.54
d-RAND-be	22.01 (11)	4.36	36.18 (19)	4.16
VCA	2.01	18.38	1.86	20.11
d-VCA	26.89 (15)	5.83	29.06 (14)	5.05
SPA	0.30	9.58	0.30	9.45
d-SPA	24.37 (13)	4.67	28.61 (14)	4.62
SNPA	24.34	9.63	36.72	5.64
d-SNPA	23.04 (13)	4.94	27.94 (13)	3.97
H2NMF	19.02	5.81	22.35	5.47
d-H2NMF	26.66 (15)	4.05	28.92 (14)	4.24
FGNSR-100	2.73	5.58	2.55	4.62
d-FGNSR-100	26.72 (14)	4.36	20.81 (8)	4.04

Table: Numerical results for the Urban data set.