

# Implicit Regularization in Regularized (Nonnegative) Low-Rank Approximations

**Jérémy Cohen**

CNRS, CREATIS

Equipe Images, LTCl, Palaiseau, March 2024

## A LASSO formulation

- ▶ Data  $y \in \mathbb{R}^m$
- ▶ Regressor  $A \in \mathbb{R}^{m \times n}$
- ▶ unknown variables  $x \in \mathbb{R}^n$

$$x_{\mu}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \|y - Ax\|_2^2 + \mu \|x\|_1 \quad (1)$$

Hyperparameter  $\mu$  controls the sparsity level of the solution  $X_{\mu}^*$ .

Figure: Polynomial fitting  $y_i = \sum_{k=0}^5 z_i^k x_k$ . Right: coefficients.

## Nice properties of LASSO for the user

- ▶ We have guarantees on the sparsity of the solution.

Let  $S$  the support of  $x_\mu^*$ . If  $n \geq m$ ,  $A[:, S]$  has full column rank.

- ▶ We can scale the values of  $\mu$  from 0 to 1.

$\mu \geq \|A^T y\|_\infty$  is equivalent to  $x_\mu^* = 0$ .

- ▶ Algorithms exist to obtain solutions for all possible  $\mu$ .

Solutions  $x_\mu^*$  are elementwise piecewise affine with respect to  $\mu$ .

Moreover:

1. The problem is convex, solvers have nice convergence guarantees.
2. We can show equivalence with constrained/Basis Pursuit versions.
3. Nonnegative LASSO has similar properties.

## The exemple of double sparse NMF

- ▶ Data matrix  $Y \in \mathbb{R}^{m_1 \times m_2}$
- ▶ decomposition rank  $r \in \mathbb{N}^*$
- ▶ unknown factors  $X_i \in \mathbb{R}_+^{n_i \times r}$ ,  $i \in \{1, 2\}$

$$\min_{X_1 \geq 0, X_2 \geq 0} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1 \quad (2)$$

Figure: Block images  $Y \approx X_1 X_2^T$ . Right:  $X_1$ ,  $X_2^T$ .

# Properties of sparse NMF

No work dedicated to the characterisation of sparse NMF solutions!

We don't know how the solution behave with  $\mu_1, \mu_2$ .

- ▶ In which space live the solutions?
- ▶ How to choose the regularization parameters?

Some related works on dictionary learning (no nonnegativity) [Georgiev 2005, Aharon 2006, Gribonval 2015, Cohen 2018] deal with identifiability. In [Cohen 2019] we show that empirically nonnegativity can help.

# Homogeneous Regularized Scale Invariant problem

We study a larger class of problems than sparse NMF, coined HRSI

$$\min_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \leq n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r g_i(X_i[:, q]) \quad (3)$$

- ▶  $f$  is a **scale-invariant** differentiable cost (e.g.  $\|Y - X_1 X_2^T\|_F^2$ ).  
Scaling  $\{X_i\}_{i \leq n}$ :  $X_i \Lambda_i$  with diagonal  $\Lambda_i$  and  $\prod_{i \leq n} \Lambda_i = I_r$ .
- ▶  $g_i$  are **homogeneous** regularization functions (e.g.  $\ell_p$  norms).
- ▶  $\mu_i$  are nonnegative regularization hyperparameters.
- ▶  $X_i$  are the unknown factors,  $X_i[:, q]$  their  $q$ -th column.

It covers sparse NMF, regularized Canonical Polyadic Decomposition, Nonnegative Tucker Decompositions...

## Takeway message

Scale invariance induces implicit penalization in HRSI.

- ▶ Better understanding of how to choose regularizations  $g_i$ , and hyperparameters  $\mu_i$ .
- ▶ Better algorithms that converge faster in loss function.
- ▶ No solution characterisation yet.

# Outline

- 1 Implicit balancing in HRSI
- 2 Explicit algorithmic balancing in alternating algorithms
- 3 Showcases on sNMF, rCPD and sNTD



# Implicit balancing in Homogeneous Regularized Scale Invariant models

## Optimizing scale in HRSI

$$\inf_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}} \phi(\{X_i\}_{i \leq n}) = \inf_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}, \prod_{i \leq n} \Lambda_i = 1} \phi(\{X_i \Lambda_i\}_{i \leq n})$$

with  $\phi$  the HRSI cost function.

- ▶ The loss in HRSI is separable with respect to scales of factors. We consider the case  $r = 1$  wlog. ( $X_i \rightarrow x_i$ )
- ▶ Since  $f$  is scale invariant, scaling the factors, i.e.

$$\prod_{i \leq n} \min_{\lambda_i = 1, \lambda_i > 0} f(\{\lambda_i x_i\}_{i \leq n}) + \sum_{i \leq n} \mu_i g_i(\lambda_i x_i) \quad (4)$$

for fixed values of  $\{x_i\}_{i \leq n}$  means finding the optimal scales to minimize the penalties

$$\prod_{i \leq n} \min_{\lambda_i = 1, \lambda_i > 0} \sum_{i \leq n} \lambda_i^{p_i} \underbrace{\mu_i g_i(x_i)}_{a_i} \quad (5)$$

where  $p_i$  is the homogeneity degree of  $g_i$ .

## A geometric mean identity

$$\min_{\forall i \leq n, \lambda_i > 0} \sum_{i \leq n} \lambda_i \text{ tel que } \prod_{i \leq n} \lambda_i = p \quad (6)$$

for  $p \geq 0$  is solved uniquely by  $\lambda_i^* = p^{1/n}$  for all  $i \leq n$ .

We can prove similarly that

$$\min_{\forall i \leq n, \lambda_i \geq 0} \sum_{i=1}^n \lambda_i^{p_i} a_i \text{ such that } \prod_{i=1}^n \lambda_i = 1 \quad (7)$$

has solutions

$$\lambda_i^* = \frac{\beta}{p_i a_i} \quad (8)$$

where  $\beta$  is the geometric mean of  $\{p_i a_i, \frac{1}{p_i}\}_{i \leq n}$

$$\beta = \left( \prod_{i \leq n} (p_i a_i)^{\frac{1}{p_i}} \right)^{\frac{1}{\sum_{i \leq n} \frac{1}{p_i}}} . \quad (9)$$

## From explicit to implicit regularization

Therefore, solutions to the scaling problem of HRSI

$$\min_{\forall i \leq n, \Lambda_i \in \mathbb{R}_+^{r \times r}} f(\{X_i \Lambda_i\}_{i \leq n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r g_i(\Lambda_i[q, q] X_i[:, q]) \quad (10)$$

where  $\Lambda_i$  are diagonal matrices such that  $\prod_{i \leq n} \Lambda_i = I_r$  are given by

$$X_i^*[:, q] = \left( \frac{\beta_q}{p_i \mu_i g_i(X_i[:, q])} \right)^{1/p_i} X_i[:, q] \quad (11)$$

Injecting this in HRSI yields an implicit formulation of HRSI:

$$\min_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \leq n}) + \tilde{\mu} \sum_{q=1}^r \left( \prod_{i=1}^n g_i(X_i[:, q])^{\frac{1}{p_i}} \right)^{\sum_{i=1}^n \frac{1}{p_i}} \quad (12)$$

with  $\tilde{\mu}$  averaged from  $\{\mu_i\}_{i \leq n}$ .

# Our main result

Proposition [C., Leplat, in preparation]

Implicit HRSI

$$\min_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \leq n}) + \tilde{\mu} \sum_{q=1}^r \left( \prod_{i=1}^n g_i(X_i[:, q])^{1/p_i} \right)^{\sum_{i=1}^n \frac{1}{p_i}} \quad (13)$$

is essentially equivalent to explicit HRSI

$$\min_{\forall i \leq n, X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \leq n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r g_i(X_i[:, q]) \quad (14)$$

- ▶ Solutions are balanced:  $\forall i \leq n, p_i \mu_i g_i(X_i[:, q]^*) = \beta_q$ .
- ▶ Implicit HRSI is fully scale-invariant.
- ▶ The nature of the regularization can change from explicit to implicit HRSI.
- ▶ Only an average  $\tilde{\mu}$  of parameters  $\mu_i$  matters!!

## Implicit regularization, ridge matrix

It is known [Srebro 2008] that

$$\operatorname{argmin}_{X_1, X_2 \in \mathbb{R}^{n_i \times r}} \|Y - X_1 X_2^T\|_F^2 + \lambda \left( \|X_1\|_F^2 + \|X_2\|_F^2 \right) \quad (15)$$

has essentially the same solutions as

$$\operatorname{argmin}_{L \in \mathbb{R}^{n_1 \times n_2}, \operatorname{rank}(L) \leq r} \|Y - L\|_F^2 + \alpha \|L\|_* \quad (16)$$

From our result, we get the implicit formulation

$$\operatorname{argmin}_{\operatorname{rank}(L_q)=1} \|M - \sum_{q \leq r} L_q\|_F^2 + \alpha \sum_{q \leq r} \|L_q\|_F. \quad (17)$$

with  $L = \sum_q L_q = X_1 X_2^T$ .

$\ell_2$  regularization in explicit HRSI induces low-rank solutions!

Mentionned in [Uschmajew 2012] for CP decomposition.

## Implicit regularization, sparse NMF

$$\min_{X_1 \geq 0, X_2 \geq 0} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1 \quad (18)$$

The implicit HRSI model for sNMF writes

$$\min_{L_q \in \mathbb{R}_+^{m_1 \times m_2}, \text{rank}(L_q) \leq 1} KL(Y, \sum_{q=1}^r L_q) + \frac{\sqrt{\mu_1 \mu_2}}{2} \sum_{q=1}^r \sqrt{\|L_q\|_1}. \quad (19)$$

- ▶ The individual sparsity levels  $\mu_i$  don't matter?!
- ▶ Sparsity occurs at the level of rank-one components.
- ▶ Using  $\ell_1^2$  removes the square-root in the implicit regularization.

### Open question

Can we use this implicit formulation to characterise solutions?

Already investigated in [Papalexakis 2013].

## Explicit balancing in alternating algorithms



## Balancing in practice

Observation: An alternating minimization algorithm can be extremely slow to converge to balanced solutions.

### Idea

Explicitly normalize factors to minimize the loss with respect to scalings.

The normalization for  $X_i$  is of the form

$$X_i[:, q]^* = \left( \frac{\beta}{p_i g_i(X_i[:, q])} \right)^{1/p_i} X_i[:, q]. \quad (20)$$

We can perform this operation for all  $i$  at the end of each outer loop of e.g. MU for sNMF.

# A meta algorithm for regularized LRA

*% Initialization:*

$$\eta \in \operatorname{argmin}_{\eta \geq 0} f(\{\eta X_i^{(0)}\}_{i \leq n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r \eta^{p_i} g_i(X_i^{(0)}[:, q])$$

$X_i^{(1)} \leftarrow \eta X_i^{(0)}$  for all  $i = 1, \dots, n$ .

**for**  $i = 1 : q$  **do**

$\{X_i^{(1)}[:, q]\}_{i \leq n}$  balanced using Equations (11)

**end for**

**for**  $k = 1 : \text{maxiter}$  **do**

*% Update of factors*  $\{X_i\}_{i \leq n}$

$X_i^{(k+1)} \leftarrow \operatorname{argmin}_{x \in \mathcal{X}} \bar{F}(x | X_i^{(k)})$  for all  $i = 1, \dots, n$ .

*% Local-optimal balancing*

**for**  $i = 1 : q$  **do**

$\{X_i^{(k+1)}[:, q]\}_{i \leq n}$  balanced using Equations (11)

**end for**

**end for**

**return**  $\{X_i[:, q]\}_{i \leq n}$

## Alternating optimization is slow?

Consider the toy problem with  $y \in \mathbb{R}$ , and  $0 < \lambda \leq y$ .

$$\min_{x_1 \in \mathbb{R}, x_2 \in \mathbb{R}} f(x_1, x_2) \text{ where } f(x_1, x_2) = (y - x_1 x_2)^2 + \lambda(x_1^2 + x_2^2). \quad (21)$$

Solutions are balanced,  $x_i^* = \sqrt{y - \lambda}$

The Alternating Least Squares algorithm

$$\begin{cases} x_1^{(k+1)} = \frac{x_2^{(k)} y}{x_2^{2(k)} + \lambda} \\ x_2^{(k+1)} = \frac{x_1^{(k+1)} y}{x_1^{2(k+1)} + \lambda} \end{cases} \quad (22)$$

is provably slow!

Proposition [C. Leplat, in preparation]

For  $k$  large enough and  $\lambda \ll y$ ,

$$\frac{x_1^{(k+1)} - \sqrt{y - \lambda}}{x_1^{(k)} - \sqrt{y - \lambda}} \approx 1 - 4 \frac{\lambda}{y}. \quad (23)$$

# Alternating optimization is slow? (2)

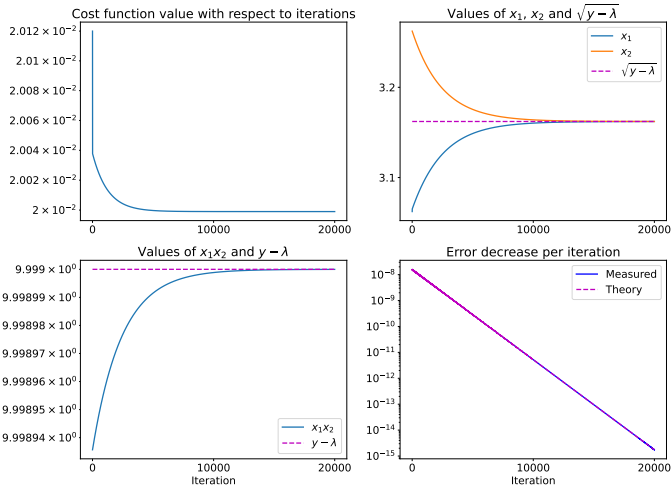


Figure:  $y = 10$  and  $\lambda = 10^{-3}$ .

**Showcases: sNMF, rCPD, sNTD**

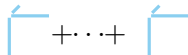
# Explicit Models



double sparse NMF (sNMF)

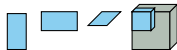
$$\min_{X_1 \in \mathbb{R}_+^{m \times r}, X_2 \in \mathbb{R}_+^{n \times r}} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1$$

ridge Canonical Polyadic Decomposition (rCPD)



$$\min_{X_i \in \mathbb{R}_+^{n+i \times r}} \|T - I_r \times_1 X_1 \times_2 X_2 \times_3 X_3\|_F^2 + \mu (\|X_1\|_F^2 + \|X_2\|_F^2 + \|X_3\|_F^2)$$

sparse Nonnegative Tucker Decomposition (sNTD)



$$\min_{\substack{W \geq 0, H \geq 0, \\ Q \geq 0, G \geq 0}} KL(T, G \times_1 W \times_2 H \times_3 Q) + \mu (\|G\|_1 + \|W\|_F^2 + \|H\|_F^2 + \|Q\|_F^2)$$

# Implicit equivalent models

## implicit sNMF

$$\min_{L_q \in \mathbb{R}_+^{m_1 \times m_2}, \text{rank}(L_q) \leq 1} KL(Y, \sum_{q=1}^r L_q) + 2\sqrt{\mu_1 \mu_2} \sum_{q=1}^r \sqrt{\|L_q\|_1}. \quad (24)$$

- ▶ Empirically, tuning  $\mu_1$  vs  $\mu_2$  has an effect, right?
- ▶ Does explicit balancing in a MU algorithm really help?

## implicit rCPD with $L_q = X_1[:, q] \otimes X_2[:, q] \otimes X_3[:, q]$

$$\min_{\{L_q\}_{1 \leq q \leq r}, \text{rank}(L_q) \leq 1} \|T - \sum_{q=1}^r L_q\|_F^2 + 3\mu \sum_{q=1}^r \|L_q\|_F^{\frac{3}{2}}. \quad (25)$$

- ▶ Empirically, do we observe a bias towards low-rank solutions?
- ▶ Does explicit balancing in a HALS algorithm really help?

## Updates for sNMF

MU factor update:

$$\hat{X}_1 = \max \left( X_1 \odot \frac{X_2 \frac{M}{X_2^T X_1}}{\mu_1 \mathbf{e}_{m_1 \times r} + \mathbf{e}_{m_1} \otimes X_2^T \mathbf{e}_{m_2}}, \epsilon \right) \quad (26)$$

$$\hat{X}_2 = \max \left( X_2 \odot \frac{\frac{M^T}{X_1 X_2^T} X_1^T}{\mu_2 \mathbf{e}_{m_2 \times r} + \mathbf{e}_{m_2} \otimes X_1 \mathbf{e}_r}, \epsilon \right) \quad (27)$$

Balancing:

$$X_1[:, q] \leftarrow \frac{\beta_q}{\mu_1 \|X_1[:, q]\|_1} X_1[:, q] \quad (28)$$

$$X_2[:, q] \leftarrow \frac{\beta_q}{\mu_2 \|X_2[:, q]\|_1} X_2[:, q] \quad (29)$$

where  $\beta_q = \sqrt{\mu_1 \mu_2 \|X_1[:, q]\|_1 \|X_2[:, q]\|_1}$ .



## Experimental Setup

XP	sNMF	rCPD
sizes $(n_i, r, \hat{r})$	$(30, 4, 4)$	$(30, 4, 6)$
data generation	$X_i \sim \mathcal{P}(\alpha X_1 X_2^T)$	$X_i \sim \mathcal{U}[0, 1]$
factors sparsity	30%	None
SNR	40	40
epsilon	$1e - 16$	$1e - 16$

sNMF: Two cases  $\mu_1 = \mu_2$  and  $\mu_1 = 1$

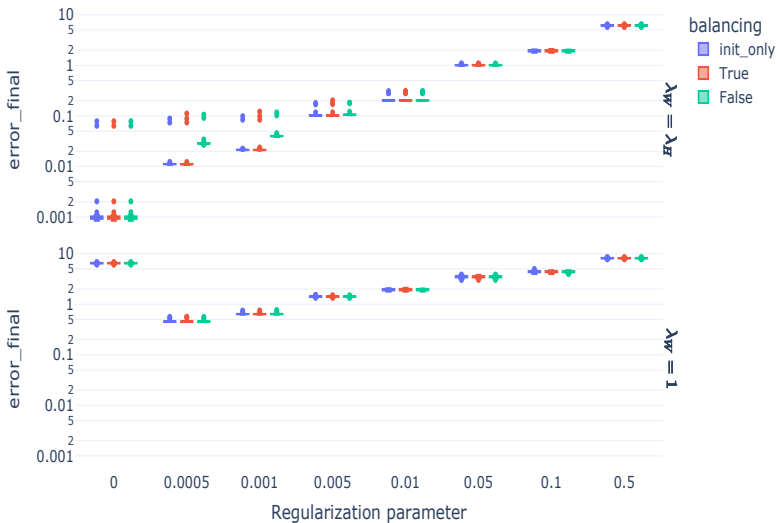
Both: Comparing balancing, no balancing and balancing at initialization only.

Evaluation with loss function, sparsity and Factor Match Score

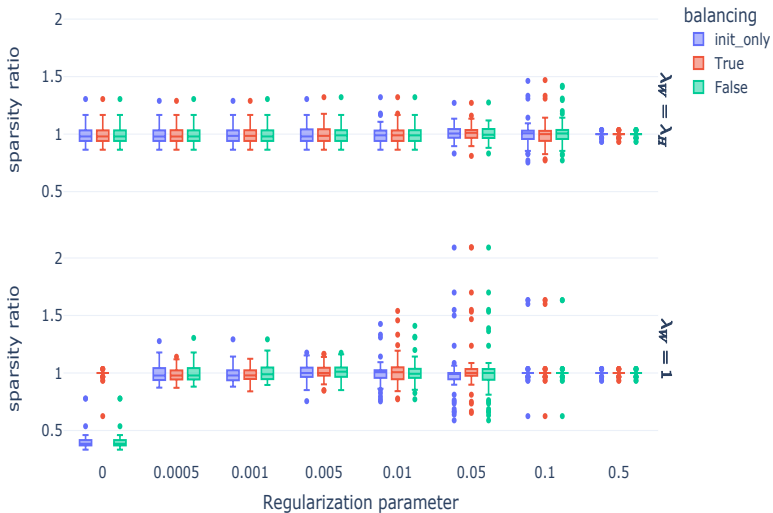
$$\text{Tr} \left( \prod_i \hat{X}_i^T X_i \right) \quad (30)$$

(after columnwise normalization and permutation).

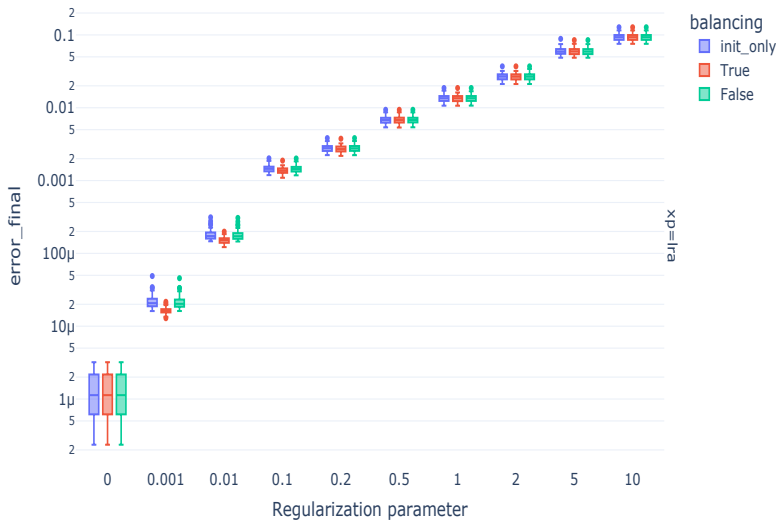
# Results for sNMF



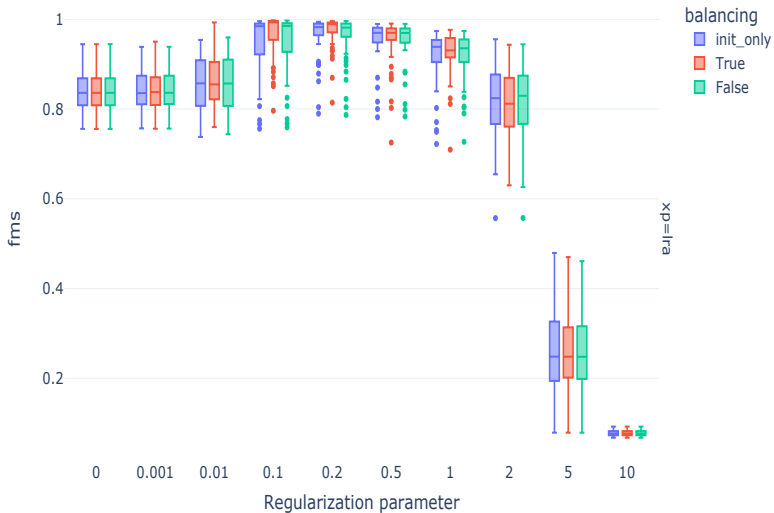
# Results for sNMF



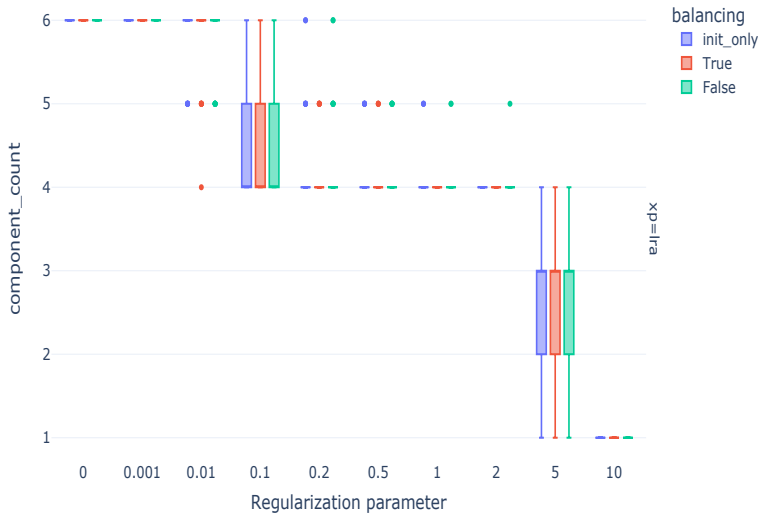
# Results for rCPD



# Results for rCPD

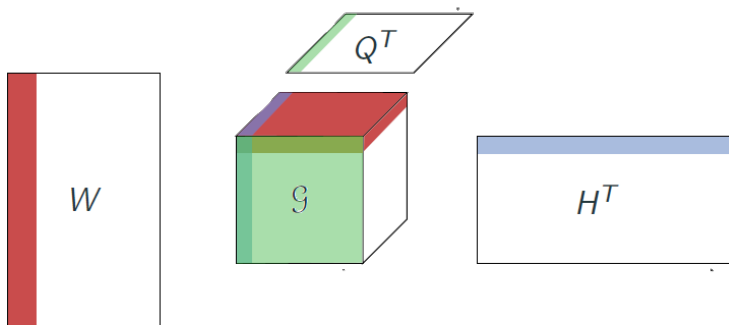


# Results for rCPD



# Sparse Nonnegative Tucker is harder

NTD does not fit HRSI because scaling ambiguity is not separable.



## Two possible scalings

### Scalar scaling

$$\min_{\substack{w \geq 0, h \geq 0, \\ q \geq 0, g \geq 0}} f(\{w, h, q, g\}) + \mu(\|g\|_1 + \|w\|_F^2 + \|h\|_F^2 + \|q\|_F^2) \quad (31)$$

where  $w = \text{vec}(W)$ ,  $h = \text{vec}(H)$ ,  $q = \text{vec}(Q)$ ,  $g = \text{vec}(G)$ .

### Sinkhorn scaling

$$\underset{\Lambda_W \text{ diagonal}}{\text{argmin}} \|\mathcal{G} \times_1 \Lambda_W \times_2 \Lambda_H \times_3 \Lambda_Q\|_1 + \|W\Lambda_W^{-1}\|_F^2 \quad (32)$$

$$\underset{\Lambda_H \text{ diagonal}}{\text{argmin}} \|\mathcal{G} \times_1 \Lambda_W \times_2 \Lambda_H \times_3 \Lambda_Q\|_1 + \|H\Lambda_H^{-1}\|_F^2 \quad (33)$$

$$\underset{\Lambda_Q \text{ diagonal}}{\text{argmin}} \|\mathcal{G} \times_1 \Lambda_W \times_2 \Lambda_H \times_3 \Lambda_Q\|_1 + \|Q\Lambda_Q^{-1}\|_F^2 . \quad (34)$$



# Ridge regularized factors estimation with KL

The problem

$$\operatorname{argmin}_{W \geq 0} KL(V|WU) + \mu \|W\|_F^2 \quad (35)$$

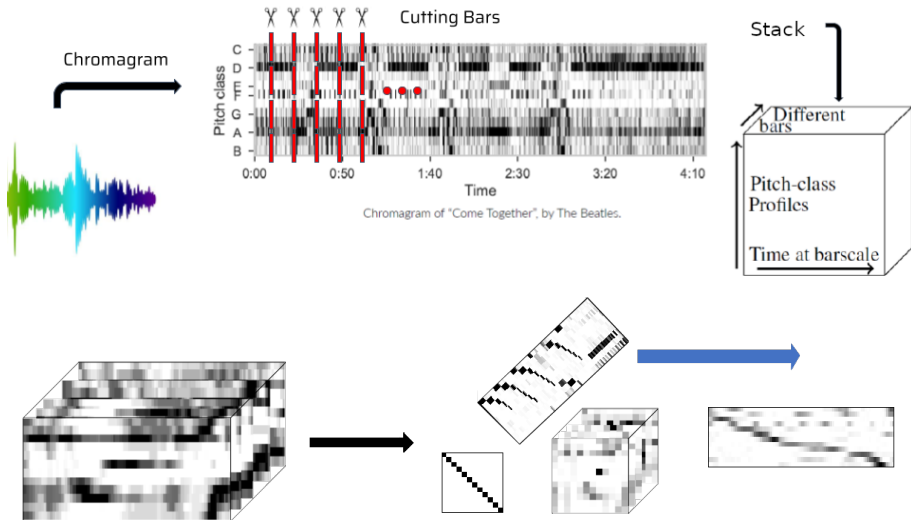
can be solved by iterating

$$\hat{W} = \max \left( \frac{[C \cdot 2 + S]^{\frac{1}{2}} - C}{2\mu}, \epsilon \right) \quad (36)$$

where  $C = EU^T$  with  $E$  is a all-one matrix of size  $m_1$ -by- $m_2 m_3$   
and  $S = 4\mu \tilde{W} \odot \left( \begin{array}{c} [V] \\ [\tilde{W}U] \end{array} U^T \right)$ .

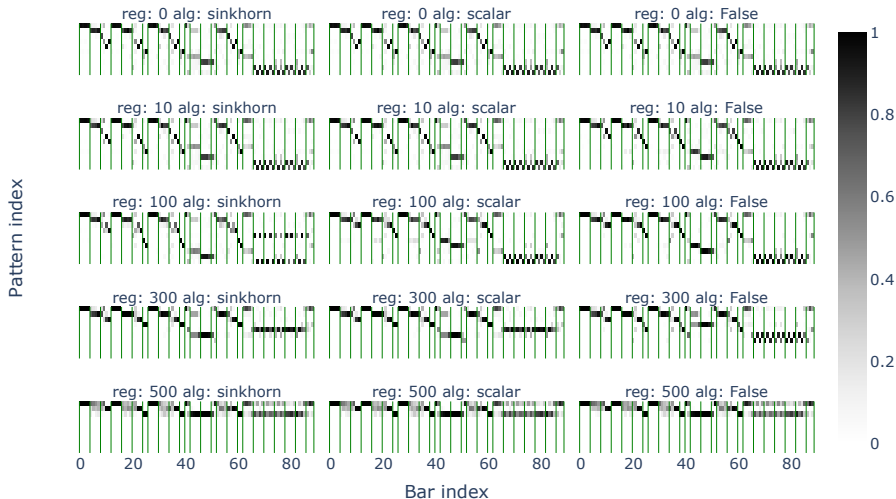
# An audio redundancy detection experiment

We perform sNTD on a tensor spectrogram, factor  $Q$  holds the song arrangement.



Note: sparsity is imposed on  $Q$  not on the core.

# An audio redundancy detection experiment



# Conclusions

## Take-aways

- ▶ We don't need normalization for regularized LRA.
- ▶ Many regularized LRA have the same solutions, don't tune the regularizations aimlessly!
- ▶ Explicit regularizations may behave unexpectedly because of scale invariance.

## Perspectives

- ▶ A clean result for linking regularized/normalized/scale-invariant HRSI?
- ▶ Optimizing the scale-invariant cost vs the explicit cost?
- ▶ A characterisation of solutions wrt sparsity levels could be obtained for rank-one?

Merci de votre attention!



# Ecole d'été Peyresq 2024

## Modèles d'approximation de rang faible et optimisation numérique

### Thématiques

- ▶ Modèles tensoriels et applications
- ▶ Relaxation et modèle d'optimisation semi-défini de rang faible
- ▶ Optimisation sur variétés et applications
- ▶ Optimisation non-convexe, algorithmes du premier ordre et unrolling
- ▶ Applications des modèles de rang faible
- ▶ Factorisation en matrices non-négatives