Implicit Regularization in Regularized (Nonnegative) Low-Rank Approximations

Jérémy Cohen, Valentin Leplat CNRS, CREATIS, Lyon

arxiv:2403.18517, under review

LORAINNE Workshop, Nancy, November 2024

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
$$X_i \in \mathbb{R}^{m_i imes r}_+$$
, $i \in \{1, 2\}$

$$\min_{X_1 \ge 0, X_2 \ge 0} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1$$
(1)

Figure: Mixture of three separable Gaussians, $\mu_1 = \mu_2 =: \mu$.

Properties of sparse NMF

No work dedicated to the characterisation of sparse NMF solutions!

We don't know how the solution behave with μ_1, μ_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

$$\min_{\forall i \le n, \ X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \le n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r g_i(X_i[:,q])$$
(2)

- ► *f* is a **scale-invariant** differentiable cost (e.g. $||Y X_1 X_2^T||_F^2$). Scaling $\{X_i\}_{i \le n}$: $X_i \Lambda_i$ with diagonal Λ_i and $\prod_{i \le n} \Lambda_i = I_r$.
- g_i are homogeneous positive-definite (e.g. l_p norms) of degree p_i.
- $\blacktriangleright \mu_i$ 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...

Takewaway message

Scale invariance induces implicit penalization in HRSI.

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

Outline



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

Our main result

Proposition [C., Leplat, in review]

Explicit HRSI

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

is essentially equivalent to an implicit HRSI problem

$$\min_{\forall i \le n, \ X_i \in \mathbb{R}^{m_i \times r}} f(\{X_i\}_{i \le n}) + \tilde{\mu} \sum_{q=1}^r \left(\prod_{i=1}^n g_i(X_i[:,q])^{\frac{1}{p_i}}\right)^{\frac{1}{\sum_{i=1}^n \frac{1}{p_i}}}$$
(4)

Also, solutions are balanced: $p_i \mu_i g_i(X_i^*[:,q])$ is constant wrt. *i*.

- 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 μ_i matters!!

Proof sketch

Solutions to the scaling problem of HRSI

$$\min_{\forall i \leq n, \Lambda_i \in \mathbb{R}^{r \times r}_+} f(\{X_i \Lambda_i\}_{i \leq n}) + \sum_{i=1}^n \mu_i \sum_{q=1}^r g_i(\Lambda_i[q, q] X_i[:, q])$$
(5)

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

$$X_{i}^{(s)}[:,q] = \underbrace{\left(\frac{\beta_{q}}{p_{i}\mu_{i}g_{i}(X_{i}[:,q])}\right)^{1/p_{i}}}_{\Lambda_{i}[q,q]^{*}}X_{i}[:,q]$$
(6)

where β_q is the geometric mean of $\{p_i \mu_i g_i(X_i)[:, q], \frac{1}{p_i}\}_{i \le n}$. Injecting this in HRSI yields an implicit formulation of HRSI:

$$\min_{X_{i}^{(s)} \in \mathbb{R}^{m_{i} \times r}} f(\{X_{i}^{(s)}\}_{i \leq n}) + \tilde{\mu} \sum_{q=1}^{r} \left(\prod_{i=1}^{n} g_{i}(X_{i}^{(s)}[:,q])^{\frac{1}{p_{i}}}\right)^{\frac{1}{\sum_{i=1}^{n} \frac{1}{p_{i}}}} \quad (7)$$

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

Implicit regularization, ridge matrix

It is known [Srebro 2008] that

$$\underset{X_{1},X_{2} \in \mathbb{R}^{n_{i} \times r}}{\operatorname{argmin}} \|Y - X_{1}X_{2}^{T}\|_{F}^{2} + \lambda \left(\|X_{1}\|_{F}^{2} + \|X_{2}\|_{F}^{2}\right)$$
(8)

has essentially the same solutions as, setting $L = X_1 X_2^T$,

$$\underset{L \in \mathbb{R}^{n_1 \times n_2}, \operatorname{rank}(L) \le r}{\operatorname{argmin}} \|Y - L\|_F^2 + \alpha \|L\|_*$$
(9)

From our result, we get the implicit formulation

$$\underset{\mathsf{rank}(L_q)=1}{\operatorname{argmin}} \|Y - L\|_F^2 + \alpha \sum_{q \le r} \|L_q\|_F.$$
(10)

with $L = \sum_{q} L_{q} = \sum_{q} X_{1}[:, q] \otimes X_{2}[:, q]^{T}$.

 ℓ_2 regularization in explicit HRSI induces low-rank solutions!

Mentionned in [Uschmajew 2012] for CP decomposition.

Implicit regularization, sparse NMF

$$\min_{X_1 \ge 0, X_2 \ge 0} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1$$
(11)

The implicit HRSI model for sNMF writes

$$\min_{L_q \in \mathbb{R}^{m_1 \times m_2}, \, \operatorname{rank}(L_q) \le 1} \, KL(Y, \sum_{q=1}^r L_q) + \frac{\sqrt{\mu_1 \mu_2}}{2} \sum_{q=1}^r \sqrt{\|L_q\|_1}.$$
(12)

- The individual sparsity levels μ_i don't matter?!
- Sparsity occurs at the level of rank-one components.

Open question

Can we use this implicit formulation to characterise solutions?

Already investigated in [Papalexakis 2013].

Explicit balancing in alternating algorithms

Alternating optimization is slow?

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

 $\min_{x_1 \in \mathbb{R}, x_2 \in \mathbb{R}} f(x_1, x_2) \text{ where } f(x_1, x_2) = (y - x_1 x_2)^2 + \lambda (x_1^2 + x_2^2) .$ (13)

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

The Alternating Least Squares algorithm

$$x_{1}^{(k+1)} = \frac{x_{2}^{(k)}y}{x_{2}^{2(k)}+\lambda}$$

$$x_{2}^{(k+1)} = \frac{x_{1}^{(k+1)}y}{x_{1}^{2(k+1)}+\lambda}$$
(14)

13/34

is provably slow!

Proposition [C. Leplat, in review] For k large enough and $\lambda \ll y$, $\frac{x_1^{(k+1)} - \sqrt{y - \lambda}}{x_1^{(k)} - \sqrt{y - \lambda}} \approx 1 - 4\frac{\lambda}{y}.$ (15)

Alternating optimization is slow? (2)



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

Balancing in practice

<u>Observation</u>: An alternating minimization algorithm can be extremely slow to converge to balanced solutions.

Idea

Explicitly normalize factors to minize the loss wrt. scalings.

The normalization for X_i is of the form

$$X_{i}[:,q]^{*} = \left(\frac{\beta}{p_{i}g_{i}(X_{i}[:,q])}\right)^{1/p_{i}}X_{i}[:,q].$$
 (16)

where β is the weighted geometric mean of the regularizations. We can perform this operation for each parameter matrix at the end of each outer loop of e.g. MU for sNMF.

A meta algorithm for regularized LRA

We propose a generic alternating Majorization Minimization algorithm with balancing, that converges in cost and iterates.

% Main optimization loop: for k = 1: maxiter do for i = 1: n and q = 1: r do $X_i[:, q]^{(k+1)} \leftarrow \operatorname{argmin}_{x \in \mathcal{X}_i} \overline{\phi}(x|X_i[:, q]^{(k)})$ end for $\{X_i^{(k+1)}[:, q]\}_{i \leq n}$ balanced using (16) for all $q \leq r$ end for return $\{X_i[:, q]\}_{i \leq n}$

Bonus contribution

We derive (closed-form) updates for ℓ_p^p regularizations with beta-divergence loss. E.g. for KL divergence and p = 2,

$$x[k] \leftarrow \frac{\sqrt{c^2 + 8\mu x[k]t} - c}{4\mu} \tag{17}$$

with $\mu > 0$ and c, t depend on data and X_i .

Showcases: sNMF, rCPD, sNTD

Explicit Models



double sparse NMF (sNMF)

$$\min_{X_1 \in \mathbb{R}^{m_1 \times r}_+, X_2 \in \mathbb{R}^{m_2 \times r}_+} KL(Y, X_1 X_2^T) + \mu_1 \|X_1\|_1 + \mu_2 \|X_2\|_1$$

ridge Canonical Polyadic Decomposition (rCPD)

$$\min_{X_{i} \in \mathbb{R}^{m_{i} \times r}_{+}} \|Y - I_{r} \times_{1} X_{1} \times_{2} X_{2} \times_{3} X_{3}\|_{F}^{2} + \mu \left(\|X_{1}\|_{F}^{2} + \|X_{2}\|_{F}^{2} + \|X_{3}\|_{F}^{2} \right)$$

sparse Nonnegative Tucker Decomposition (sNTD) in the paper!

Implicit equivalent models

implicit sNMF

$$\min_{L_q \in \mathbb{R}^{m_1 \times m_2}, \, \operatorname{rank}(L_q) \le 1} \, KL(Y, \sum_{q=1}^r L_q) + 2\sqrt{\mu_1 \mu_2} \sum_{q=1}^r \sqrt{\|L_q\|_1}.$$
(18)

Empirically, tuning µ₁ vs µ₂ has an effect, right?
 Does explicit balancing in a MU algorithm really help?

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

$$\min_{\{L_q\}_{1\leq r}, \text{ rank}(L_q)\leq 1} \|Y - \sum_{q=1}^r L_q\|_F^2 + 3\mu \sum_{q=1}^r \|L_q\|_F^{\frac{2}{3}}.$$
 (19)

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

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

<u>sNMF</u>: Two cases $\mu_1 = \mu_2$ and $\mu_1 = 1$

<u>Both</u>: Comparing balancing, no balancing and balancing at initialization only.

Evaluation with loss function, sparsity and Factor Match Score

$$\operatorname{Tr}\left(\prod_{i} \hat{X}_{i}^{T} X_{i}\right) \tag{20}$$

(after columnwise normalization and permutation).

Results for sNMF



Results for sNMF



Results for rCPD



Results for rCPD



Conclusions

Take-aways

- Explicit regularizations may behave unexpectedly because of scale invariance.
- Don't tune the hyperparameters independently!
- Balance your solutions in the algorithm, or at least at first and last iteration.

Perspectives

What about non-homogeneous/non positive definite regularizations?

Thank you for your attention! 🖌



Updates for sNMF

MU factor update:

$$\hat{X}_{1} = \max\left(X_{1} \odot \frac{X_{2} \frac{M}{X_{2}^{T} X_{1}}}{\mu_{1} e_{m_{1} \times r} + e_{m_{1}} \otimes X_{2}^{T} e_{m_{2}}}, \epsilon\right)$$
(21)
$$\hat{X}_{2} = \max\left(X_{2} \odot \frac{\frac{M^{T}}{X_{1} X_{2}^{T}} X_{1}^{T}}{\mu_{2} e_{m_{2} \times r} + e_{m_{2}} \otimes X_{1} e_{r}}, \epsilon\right)$$
(22)

Balancing:

$$X_{1}[:,q] \leftarrow \frac{\beta_{q}}{\mu_{1} \|X_{1}[:,q]\|_{1}} X_{1}[:,q]$$
(23)
$$X_{2}[:,q] \leftarrow \frac{\beta_{q}}{\mu_{2} \|X_{2}[:,q]\|_{1}} X_{2}[:,q]$$
(24)

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

Results for rCPD



Sparse Nonnegative Tucker is harder

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



Two possible scalings

Scalar scaling

$$\min_{\substack{w \ge 0, h \ge 0, \\ q \ge 0, g \ge 0}} f(\{w, h, q, g\}) + \mu(\|g\|_1 + \|w\|_F^2 + \|h\|_F^2 + \|q\|_F^2)$$
(25)

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

Sinkhorn scaling

Ridge regularized factors estimation with KL

The problem

$$\underset{W \ge 0}{\operatorname{argmin}} KL(V|WU) + \mu \|W\|_F^2$$
(29)

can be solved by iterating

$$\hat{W} = \max\left(\frac{\left[C^{.2} + S\right]^{.\frac{1}{2}} - C}{2\mu}, \epsilon\right)$$
(30)

where $C = EU^T$ with E is a all-one matrix of size m_1 -by- m_2m_3 and $S = 4\mu \tilde{W} \odot \left(\frac{[V]}{[\tilde{W}U]}U^T\right)$. An audio redundancy detection experiment We perform sNTD on a tensor spectrogramm, factor Q holds the song arrangement.



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

An audio redundancy detection experiment



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 ϕ 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 → x_i)
- Since *f* is scale invariant, scaling the factors, i.e.

$$\min_{\prod_{i\leq n}\lambda_i=1,\ \lambda_i>0} f(\{\lambda_i x_i\}_{i\leq n}) + \sum_{i\leq n} \mu_i g_i(\lambda_i x_i)$$
(31)

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

$$\min_{\prod_{i\leq n}\lambda_i=1,\ \lambda_i>0} \sum_{i\leq n} \lambda_i^{p_i} \underbrace{\mu_i g_i(x_i)}_{a_i}$$
(32)

where p_i is the homogeneity degree of g_i .

A geometric mean identity

$$\min_{\forall i \le n, \ \lambda_i > 0} \ \sum_{i \le n} \lambda_i \text{ tel que } \prod_{i \le n} \lambda_i = p$$
(33)
for $p \ge 0$ is solved uniquely by $\lambda_i^* = p^{1/n}$ for all $i \le n$.

We can prove similarly that

$$\min_{\forall i \le n, \ \lambda_i \ge 0} \sum_{i=1}^n \lambda_i^{p_i} a_i \text{ such that } \prod_{i=1}^n \lambda_i = 1$$
(34)

has solutions

$$\lambda_i^* = \frac{\beta}{p_i a_i} \tag{35}$$

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

$$\beta = \left(\prod_{i \le n} (p_i a_i)^{\frac{1}{p_i}}\right)^{\frac{1}{\sum_{i \le n} \frac{1}{p_i}}}.$$
(36)

34/34