Extrapolation on BCD for constrained Matrix/Tensor Factorization

- ► Andersen Ang (Postdoc in Dept. Combinatorics and Optimization, UWaterloo, Canada)
- Joint work with
 - ► Jérémy E. Cohen (CNRS, IRISA Rennes, France)
 - ► Nicolas Gillis, Le Thi Khanh Hien (UMONS, Belgium)
 - ► Valentin Leplat (UCLouvain, Belgium)
- 1. (On Complex NMF) **A.**, Leplat, Gillis, "Fast algorithm for complex-NMF with application to source separation", Submitted to conference EUSIPCO21, February 2021
- (On NTF/NCPD) A., Cohen, Gillis, Hien, "Accelerating Block Coordinate Descent for Nonnegative Tensor Factorization", Journal Numerical Linear Algebra with Applications, 2021
- 3. (On CPD) **A.**, Cohen, Hien and Gillis, "Extrapolated Alternating Algorithms for Approximate Canonical Polyadic Decomposition", Conference IEEE ICASSP 2020
- 4. (On NMF) **A.** and Gillis, "Accelerating Nonnegative Matrix Factorization Algorithms using Extrapolation", Journal Neural Computation 31 (2), 2019

SIAM CSE21 2021-March-4

One sentence summary / take-home

An heuristic extrapolated BCD algorithm with reduced restart cost, that works fast on X-factorization.

What problem to solve

 (Nonnegative) Canonical Polyadic Decomposition¹ (CPD) with known factorization rank

• Given $r \in \mathbb{N}_+$ and $\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, find $\mathbf{A}_i \in \mathbb{R}^{I_i \times r}$ for $i = 1, 2 \dots, N$ by solving

$$(\mathcal{P}) : \min_{\mathbf{A}_1, \dots, \mathbf{A}_N} \left\| \mathcal{T} - \mathcal{I} \times_1 \mathbf{A}_1 \times_2 \mathbf{A}_2 \times_3 \dots \times_N \mathbf{A}_N \right\|_F^2$$

subject to $\mathbf{A}_i \in \mathcal{C}_i$ (for example, nonnegativity).

► Statement: try to accelerate (exact and inexact) BCD algorithm for solving (P).

¹The acceleration also works for Tucker model and even in complex numbers.

What problem to solve - simplified case

• Given r and $\mathcal{T} \in \mathbb{R}^{I \times J \times K}$, find $\mathbf{U} \in \mathbb{R}^{I \times r}$, $\mathbf{V} \in \mathbb{R}^{J \times r}$, $\mathbf{W} \in \mathbb{R}^{K \times r}$ by

$$(\mathcal{P}): \min_{\mathbf{U} \ge \mathbf{0}, \mathbf{V} \ge \mathbf{0}, \mathbf{W} \ge \mathbf{0}} f(\mathbf{U}, \mathbf{V}, \mathbf{W}) \coloneqq \left\| \mathcal{T} - \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \right\|_F^2$$

subject to, says $\mathbf{U}, \mathbf{V}, \mathbf{W}$ nonnegative.

The BCD that you know very well

Algorithm 1: A typical BCD to solve 3rd-order NCPD

Result: $\mathbf{U}, \mathbf{V}, \mathbf{W}$ that minimize f

- 1 Initialize $\mathbf{U}_0, \mathbf{V}_0, \mathbf{W}_0;$
- 2 while not converge do

$$\mathbf{U}_k = \operatorname*{argmin}_{\mathbf{U} > \mathbf{0}} f(\mathbf{U}, \mathbf{V}_{k-1}, \mathbf{W}_{k-1});$$

4
$$\mathbf{V}_k = \operatorname*{argmin}_{\mathbf{V} \ge \mathbf{0}} f(\mathbf{U}_k, \mathbf{V}, \mathbf{W}_{k-1});$$

$$\mathbf{W}_k = \underset{\mathbf{W} \ge \mathbf{0}}{\operatorname{argmin}} f(\mathbf{U}_k, \mathbf{V}_k, \mathbf{W});$$

6 end

5

- Can use what ever solver on each sub-problem: (acc-)gradient descent, active set, ADMM, (acc-)HALS, MU
- Note: exact BCD vs inexact BCD.

Algorithm 2: Heuristic Extrapolated BCD with Restarts (HER)

Result: $\mathbf{U}, \mathbf{V}, \mathbf{W}$ that minimize fInitialize $U_0, V_0, W_0, \hat{U}_0, \hat{V}_0, \hat{W}_0, \beta_0 \in (0, 1], \eta \ge \bar{\gamma} \ge \gamma \ge 1;$ 1 2 while not converge do $\mathbf{U}_{k} = \operatorname{argmin} f(\mathbf{U}, \hat{\mathbf{V}}_{k-1}, \hat{\mathbf{W}}_{k-1});$ 3 $\hat{\mathbf{U}}_{k} = [\mathbf{U}_{k} + \beta_{k-1}(\mathbf{U}_{k} - \mathbf{U}_{k-1})]_{+}$ 4 $\mathbf{V}_k = \operatorname{argmin} f(\hat{\mathbf{U}}_k, \mathbf{V}, \hat{\mathbf{W}}_{k-1});$ 5 $\hat{\mathbf{V}}_{k} = [\mathbf{V}_{k}^{\mathsf{V}} + \beta_{k-1}(\mathbf{V}_{k} - \mathbf{V}_{k-1})]_{\perp}:$ 6 $\mathbf{W}_{k} = \operatorname{argmin} f(\hat{\mathbf{U}}_{k}, \hat{\mathbf{V}}_{k}, \mathbf{W});$ 7 $\hat{\mathbf{W}}_{h} = [\mathbf{W}_{h}^{\mathsf{W}} + \beta_{h-1}(\mathbf{W}_{h} - \mathbf{W}_{h-1})]_{+}$ 8 Compute $\hat{F}_k := f(\hat{\mathbf{U}}_k, \hat{\mathbf{V}}_k, \mathbf{W}_k);$ 9 if $\hat{F}_k > \hat{F}_{k-1}$ then 10 $\hat{\mathbf{U}}_k = \mathbf{U}_k, \hat{\mathbf{V}}_k = \mathbf{V}_k, \hat{\mathbf{W}}_k = \mathbf{W}_k;$ 11 $\bar{\beta}_k = \beta_{k-1}, \beta_k = \beta_{k-1}/\eta;$ 12 13 else $\mathbf{U}_k = \hat{\mathbf{U}}_k, \mathbf{V}_k = \hat{\mathbf{V}}_k, \mathbf{W}_k = \hat{\mathbf{W}}_k;$ 14 $\beta_k = \min\{\bar{\beta}_k, \beta_{k-1}\gamma\}, \ \bar{\beta}_k = \min\{1, \bar{\beta}_{k-1}\bar{\gamma}\};$ 15 16 end 17 end



Figure: Average over 100 run of ANLS, A-HALS and their extrapolated variants applied on low-rank (left) synthetic data sets, the method of $APG-MF^2$ is used for comparisons.

Similar figure for many other cases ...

 2 Xu, Y., Yin, W.: A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. SIAM J Imaging Sciences 6(3), 2013 6 / 15

It works well in NTF: CPD form

$$\min_{\mathbf{U} \ge \mathbf{0}, \mathbf{V} \ge \mathbf{0}, \mathbf{W} \ge \mathbf{0}} f(\mathbf{U}, \mathbf{V}, \mathbf{W}) \coloneqq \left\| \mathcal{T} - \mathbf{U} \times_2 \mathbf{V} \times_3 \mathbf{W} \right\|_F^2$$



Figure: On 50^3 tensor, rank 10. Blue : ordinary BCD. Orange : HER-BCD.

Similar figure for AS, accPGD, GD, MU ...

It works well in NTF: other cases

Fig.	Test description	$[I_1, I_2, I_3, r, \sigma]$
Synthetic data		
2	Cube size, low rank, noiseless	[50, 50, 50, 10, 0]
2	Unbalanced size, low rank, noiseless	$\left[150, 10^3, 50, 12, 0 ight]$
3	Unbalanced size, larger rank, noiseless	$\left[150, 10^3, 50, 25, 0 ight]$
4	Large cube size, low rank, noisy	$\left[500, 500, 500, 10, 0.01 ight]$
5	Unbalanced size, low rank, noisy, ill-condition	$[150, 10^3, 50, 12, 0.001]$
6	HER-AO-gradients compared with APG and iBPG	$[150, 10^3, 50, 10, 0.01]$
7	Comparing {HER,Bro,GR,LS}-AHALS	$[50, 50, 50, 10, 0] \ [150, 10^3, 50, 12, 0.01] \ [150, 10^3, 50, 25, 0.01]$
Real data		
8	Two HSI images : PaviaU and Indian Pine	$\begin{matrix} [610, 340, 103, 10] \\ [145, 145, 200, 15] \end{matrix}$
9	Big data : black-and-white video sequence	$[153, 238, 1.4 \times 10^4, \{10, 20, 30\}]$
Table 4: List of experiments on NTF.		

Same result for AS, accPGD, GD, MU ... So many test cases because NO theory.



Figure: On a $153\times238\times14000$ video tensor that has very high mode-wise condition numbers, rank 30. 9 / 15

It works well in CPD



Figure: On Wine data $44 \times 2700 \times 200$, rank 15.

Note: here the problem has no constraint, block-wise subproblem has closed-form solution (ALS, which is also Newton's iteration), so this shows that HER also works with second-order method. 10/15



Figure: On fitting a complex matrix (spectrogram of speech data), rank 75.

Algorithm 2: HER

- 1: Input: a nonnegative N-way tensor
- 2: Output: nonnegative factors $A^{(1)}, A^{(2)}, \ldots, A^{(N)}$.
- 3: Initialization: Choose $\beta_0 \in (0, 1), \eta \geq \bar{\gamma} \geq \gamma \geq 1$ and 2 sets of initial factor matrices $(A_0^{(1)}, \ldots, A_0^{(N)})$ and $(\hat{A}_0^{(1)}, \ldots, \hat{A}_0^{(N)})$. Set $\bar{\beta}_0 = 1$ and k = 1.

4: repeat

- 5: for $i = 1, \ldots, N$ do
- 6: **Update step** Let $A_k^{(i)}$ be an exact/inexact solution of

$$\min_{\mathbf{A}^{(i)} \ge 0} F\left(\hat{A}_{k}^{(1)}, \dots, \hat{A}_{k}^{(i-1)}, A^{(i)}, \hat{A}_{k-1}^{(i+1)}, \dots, \hat{A}_{k-1}^{(N)}\right).$$
(15)

7: Extrapolation step

$$\hat{A}_{k}^{(i)} = \max\left(0, A_{k}^{(i)} + \beta_{k-1}(A_{k}^{(i)} - A_{k-1}^{(i)})\right).$$
(16)

8: end for

9: Compute $\hat{F}_k := F\left(\hat{A}_k^{(1)}, \hat{A}_k^{(2)}, \dots, \hat{A}_k^{(N-1)}, A_k^{(N)}\right).$

10: if
$$F_k > F_{k-1}$$
 then

11: Set
$$\hat{A}_{k}^{(i)} = A_{k}^{(i)}$$
, $i = 1, ..., N$ % abandon the sequence $\hat{A}_{k}^{(i)}$

12: Set
$$\beta_k = \beta_{k-1}$$
, $\beta_k = \beta_{k-1}/\eta$. % Update β , decrease

14: Set
$$A_k^{(i)} = \hat{A}_k^{(i)}, i = 1, ..., N. \%$$
 keep the sequence $\hat{A}_k^{(i)}$

15: Set
$$\beta_k = \min\{1, \beta_{k-1}\bar{\gamma}\}, \quad \beta_k = \min\{\beta_{k-1}, \beta_{k-1}\gamma\}.$$
 % Increase β and β

- 16: end if
- 17: Set k = k + 1.
- 18: **until** some criteria is satisfied



13/15

What are the important points

- ▶ Why no theory: it is accelerating nonconvex BCD, not Gradient Descent, not many useful tools to use, so cannot prove convergence
- Why faster: acceleration due to extrapolation with restarts (safe-guard mechanism) that is cheap to compute.
- Why restart cost is low: the use of *F̂* reuse already computed component when updating the blocks The reduction: from *O*(∏^N_{i=1} *I_i*) to *O*(*r^{N-1}I_N*)
- ▶ Use of \hat{F} as a surrogate of F makes sense: $\lim_{k \to \infty} |\hat{F}_k F_k| = 0.$
- What about other findings
 - ▶ BCD can be replaced by inexact BCD, even 1 step gradient descent
 - Constrained problem can be replaced by unconstrained problem
 - Tensor can be replaced by matrix
 - Even works for complex number



The End. Papers, code, slide available at angms.science