

# CPD with Structural Constraints: From BCD to Stochastic Mirror Descent

Xiao Fu

<sup>‡</sup>School of EECS, Oregon State University

**SIAM LA 2021, Virtual Conference**

ack: Shahana Ibrahim, Kejun Huang, Wenqiang Pu, Mingyi Hong, Hoi-To Wai, Nico Vervilet,  
Nicolas Gillis and Lieven De Lathauwer

August 13, 2021

This talk is based on

- ▶ X. Fu, N. Vervliet, L. De Lathauwer, K. Huang and N. Gillis, "[Computing Large-Scale Matrix and Tensor Decomposition With Structured Factors: A Unified Nonconvex Optimization Perspective](#)," in IEEE Signal Processing Magazine, vol. 37, no. 5, pp. 78-94, Sept. 2020.
- ▶ X. Fu, S. Ibrahim, H. Wai, C. Gao and K. Huang, "[Block-Randomized Stochastic Proximal Gradient for Low-Rank Tensor Factorization](#)," in IEEE Transactions on Signal Processing, vol. 68, pp. 2170-2185, 2020
- ▶ W. Pu, S. Ibrahim, X. Fu, and M. Hong, "[Stochastic Mirror Descent for Low-Rank Tensor Decomposition Under Non-Euclidean Losses](#)", submitted to IEEE Transaction on Signal Processing, April, 2021.



# Canonical Polyadic Decomposition (CPD)

$\mathcal{T}$

$\approx$

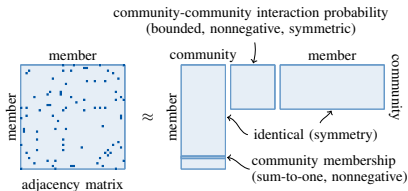
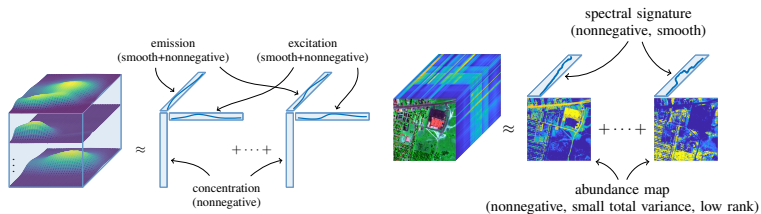
$A_1(:, 1)$   $A_1(:, 2)$   $A_1(:, 3)$

$+ \dots +$

$A_1(:, R)$   $A_2(:, R)$   $A_3(:, R)$

$\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$   $A_n \in \mathbb{R}^{I_n \times R}, n = 1, 2, 3$

# Decomp. with Constraints/Regularization



# CPD with Constraints/Regularization

General problem of interest:

$$\begin{array}{ll} \min_{\text{model param.}} & \text{dist}(\text{data}||\text{model}) + \left( \begin{array}{l} \text{penalty for} \\ \text{structure violation} \end{array} \right) \\ \text{under} & \text{structural constraints,} \end{array}$$

For example [Beutel et al., 2014]:

$$\begin{array}{ll} \min_{\{\mathbf{A}_n\}_{n=1}^N} & \frac{1}{2} \|\mathcal{T} - \llbracket \mathbf{A}_1, \dots, \mathbf{A}_N \rrbracket\|_{\text{F}}^2 + \lambda \sum_{n=1}^N \|\mathbf{A}_n\|_1 \\ \text{s.t.} & \mathbf{A}_n \geq 0. \end{array}$$

# CPD with Constraints/Regularization

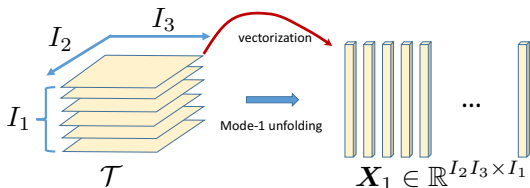
Structural info. on latent factors are useful for

- ▶ fending against noise (acting as priors);
- ▶ enhancing interpretability;
- ▶ avoiding ill-posedness [Lim and Comon, 2009];
- ▶ improving identifiability (esp.in matrix fact.) [Fu et al., 2019a];

Frequently seen constraints/regularization:

- ▶ nonnegativity:  $\mathbf{A}_n \geq 0$  [Chi and Kolda, 2012]
- ▶ sparsity:  $\|\mathbf{A}_n\|_1$
- ▶ prob. simplex:  $\mathbf{1}^\top \mathbf{A}_n = \mathbf{1}^\top$ ,  $\mathbf{A}_n \geq 0$  [Kargas et al., 2018, Yeredor and Haardt, 2019]
- ▶ boundedness:  $a \leq \mathbf{A}_n(i, r) \leq b$
- ▶ column/row sparsity  $\|\mathbf{A}_n\|_{2,1}$  or  $\|\mathbf{A}_n\|_{1,2}$  [Yang et al., 2015].
- ▶ **more**: monotonicity, smoothness, total variation, orthogonality, symmetry ... see [Sidiropoulos et al., 2017]

## First glance: Not so hard?



For all  $i_1, \dots, i_N$ , define *matrix unfolding*:

$$\mathbf{X}_n(j, i_n) = \mathcal{T}(i_1, \dots, i_N),$$

where  $j = 1 + \sum_{\ell=1, \ell \neq n} (i_\ell - 1)J_\ell$ ,  $J_\ell = \prod_{m=1, m \neq n}^{\ell-1} I_m$ .

$$\mathbf{X}_n = \mathbf{H}_n \mathbf{A}_n^\top,$$

where the matrix  $\mathbf{H}_n \in \mathbb{R}^{(\prod_{\ell=1, \ell \neq n}^N I_\ell) \times R}$  is defined as:

$$\mathbf{H}_n = \mathbf{A}_N \odot \dots \odot \mathbf{A}_{n+1} \odot \mathbf{A}_{n-1} \odot \dots \odot \mathbf{A}_1.$$

# First glance: Not so hard?

BCD updates:

$$\mathbf{A}_n^{(t+1)} \leftarrow \arg \min_{\mathbf{A}_n} \frac{1}{2} \|\mathbf{X}_n - \mathbf{H}_n^{(t)} \mathbf{A}^\top\|_F^2 + h_n(\mathbf{A}_n), \quad (1)$$

where  $\mathbf{H}_n^{(t)} = \mathbf{A}_N^{(t)} \odot \dots \odot \mathbf{A}_{n+1}^{(t)} \odot \mathbf{A}_{n-1}^{(t+1)} \odot \dots \odot \mathbf{A}_1^{(t+1)}$ , since  $\mathbf{A}_\ell$  for  $\ell < n$  has been updated.

- ▶ The subproblem is often convex - “easy” to solve.
  - ▶ ADMM [[Huang et al., 2016](#)]
  - ▶ proximal/projected gradient descent [[Lin, 2007](#)]
  - ▶ accelerated proximal/projected gradient descent [[Liavas et al., 2017](#)]
  - ▶ inexact accelerated gradient descent [[Xu and Yin, 2013](#)]
- ▶ Convergence well understood [[Razaviyayn et al., 2013](#)].



# Bottleneck: MTTKRP

**Common challenge:** matricized tensor times Khatri–Rao product (MTTKRP), i.e.,  $\mathbf{X}_n^\top \mathbf{H}_n^{(t)}$ .

- ▶ This single step costs  $O(\prod_{n=1}^N I_n R)$  flop.
- ▶ It also could use up to  $O(\prod_{n=1}^N I_n)$  memory (depending on implementation).

## Memory-efficient implementations:

- ▶ Large-scale sparse tensor [[Phipps and Kolda, 2019](#), [Smith et al., 2015](#)]
- ▶ Large-scale dense tensor [[Ravindran et al., 2014](#), [Kolda and Bader, 2006](#)]
- ▶ `tensor_toolbox` has nice modules for MTTKRP if you use Matlab

# Stochastic Optimization in One Slide

General stochastic optimization under finite sum:

$$\min_{\boldsymbol{\theta}} \frac{1}{L} \sum_{\ell=1}^L f_{\ell}(\boldsymbol{\theta}) + h(\boldsymbol{\theta}),$$

Stochastic proximal gradient

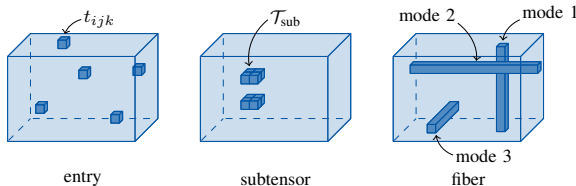
$$\boldsymbol{\theta}^{(t+1)} \leftarrow \text{Prox}_h \left( \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \mathbf{g}(\boldsymbol{\theta}^{(t)}) \right),$$

where  $\mathbf{g}(\boldsymbol{\theta}^{(t)})$  is a “stochastic oracle” evaluated at  $\boldsymbol{\theta}^{(t)}$ :

- ▶ **Example:**  $\mathbf{g}(\boldsymbol{\theta}^{(t)}) = 1/|\mathcal{S}^{(t)}| \sum_{\ell \in \mathcal{S}^{(t)}} \nabla f_{\ell}(\boldsymbol{\theta}^{(t)})$ .
- ▶  $\mathcal{S}^{(t)} \subset [L]$  is a random subset.

# Sampling Tensor Data for Stochastic Opt.

A natural way to circumvent MTTKRP - working with sampled data



- ▶ Entry sampling [Beutel et al., 2014, Hong et al., 2020, Kolda and Hong, 2020]
- ▶ Subtensor sampling [Vervliet and De Lathauwer, 2016]
- ▶ Fiber sampling [Battaglini et al., 2018]

- ▶ The Euclidean loss based CPD problem can be written as:

$$\min_{\{\mathbf{A}_n\}_{n=1}^N} \frac{1}{L} \sum_{i_1=1}^{l_1} \cdots \sum_{i_N=1}^{l_N} \underbrace{\left( \mathcal{T}(i_1, \dots, i_N) - \sum_{r=1}^R \prod_{n=1}^N \mathbf{A}_n(i_n, r) \right)^2}_{f_{i_1, \dots, i_N}(\boldsymbol{\theta})} + h(\boldsymbol{\theta})$$

where  $L = \prod_{n=1}^N l_n$ ,  $h(\boldsymbol{\theta}) = \sum_{n=1}^N h_n(\mathbf{A}_n)$  and  $\boldsymbol{\theta} = [\text{vec}(\mathbf{A}_1)^\top, \dots, \text{vec}(\mathbf{A}_N)^\top]^\top$ .

- ▶ The corresponding SGD update:

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \text{Prox}_h \left( \boldsymbol{\theta}^{(t)} - \frac{\alpha^{(t)}}{|\mathcal{B}^{(t)}|} \sum_{(i_1, \dots, i_N) \in \mathcal{B}^{(t)}} \nabla f_{i_1, \dots, i_N}(\boldsymbol{\theta}^{(t)}) \right).$$

- ▶ very small per-iteration complexity.

# Stochastic Optimization for CPD

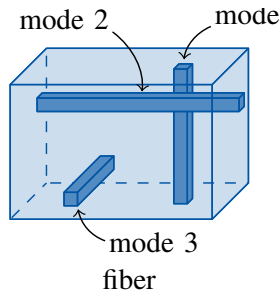
- ▶ **Challenge 1:** constraint enforcing - random sample may create some problems.
  - ▶  $\mathcal{T}(i, j, k)$  only contains info of  $\mathbf{A}_1(i, :)$ ; hard to enforce constraints like  $\mathbf{1}^\top \mathbf{A}_1(:, r) = 1$ .
- ▶ **Challenge 2:** step size scheduling - what is the best practice?

*"One of the major issues in stochastic gradient descent (SGD) methods is how to choose an appropriate step size while running the algorithm."* [Tan et al., 2016]

*"Determining a good learning rate becomes more of an art than science for many problems."* [Zeiler, 2012]
- ▶ **Challenge 3:** convergence analysis

# Fiber Sampling and Sketched LS

$$\underbrace{\mathcal{T}(i_1, \dots, i_{n-1}, :, i_{n+1}, \dots, i_N)}_{\text{a mode-}n \text{ fiber}} = \mathbf{X}_n(j_n, :) = \mathbf{H}_n(j_n, :)\mathbf{A}_n^\top$$



Sample a set of mode- $n$  fibers, indexed by  $\mathcal{Q}_n^{(t)}$  and solve a 'sketched least squares' problem [Battaglino et al., 2018]:

$$\mathbf{A}_n^{(t+1)} \leftarrow \arg \min_{\mathbf{A}_n} \left\| \mathbf{X}_n(\mathcal{Q}_n^{(t)}, :) - \mathbf{H}_n^{(t)}(\mathcal{Q}_n^{(t)}, :)\mathbf{A}_n^\top \right\|_F^2,$$

# Challenges

## Pros:

- ▶  $\mathbf{A}_n^{(t+1)} \leftarrow (\mathbf{H}_n^{(t)}(\mathcal{Q}_n^{(t)}, :)^\dagger \mathbf{X}_n(\mathcal{Q}_n^{(t)}, :))^\top$  updates the entire  $\mathbf{A}_n$ .
- ▶ No step size selection.

## Challenges:

- ▶  $|\mathcal{Q}_n^{(t)}| \geq R$  (suggested as  $|\mathcal{Q}_n^{(t)}| = 10R \log R$  in [Battaglino et al., 2018]); can be costly when  $R$  is large ( $R$  could reach  $O(I^2)$ ).
- ▶ Convergence unknown.

## Proposed Approach

For constraints/reg.: Use SGD instead of LS:

- ▶ Construct  $\mathbf{G}_n^{(t)} = \mathbf{A}_n \mathbf{B}^\top \mathbf{B} - \mathbf{X}_n(\mathcal{Q}_n, :)^{\top} \mathbf{B}$  with  $\mathbf{B} = \mathbf{H}^{(t)}(\mathcal{Q}_n^{(t)}, :)$ . Then

$$\mathbf{A}_n^{(t+1)} \leftarrow \text{Prox}_{h_n} \left( \mathbf{A}_n^{(t)} - \alpha^{(t)} \mathbf{G}_n^{(t)} \right).$$

- ▶ can deal with a large number of constraints with closed-form/semi-algebraic solutions.
- ▶ small memory footprint and lightweight in terms of flops.

### Remaining Challenges:

- ▶ Step size is back - tuning can be irritating.
- ▶  $\mathbb{E}[\mathbf{g}] \neq \nabla_{\theta} f(\theta)$ : not desired for SGD convergence, where  $\mathbf{g} = [\text{vec}(\mathbf{G}_1)^{\top}, \dots, \text{vec}(\mathbf{G}_N)^{\top}]^{\top}$ .



## Simple Fix

- ▶ **Unbiased gradient estimation** - block randomization
  - ▶ for each iteration, sample a block  $n$  to update (uniformly).
  - ▶ for the sampled block, sample corresponding fibers and do stochastic proximal gradient.
  - ▶  $\mathbb{E}[\mathbf{g}] = c\nabla_{\theta}f(\theta)$  now holds for  $c > 0$ .
- ▶ **Step size rule** - ideas from deep learning

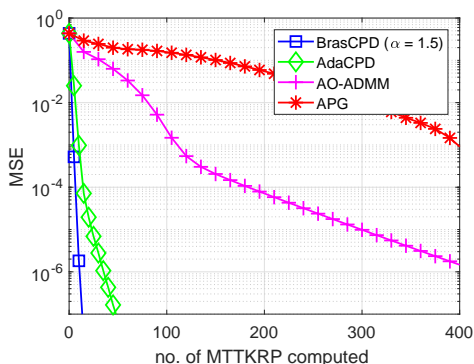
$$[\eta_n^{(t)}]_{i,r} \leftarrow \frac{1}{\left(b + \sum_{q=1}^t [\mathbf{G}_n^{(q)}]_{i,r}^2\right)^{1/2+\epsilon}},$$

where  $b$  and  $\epsilon$  are for regularization purpose.

$$\mathbf{A}_n^{(t+1)} \leftarrow \text{Prox}_{h_n} \left( \mathbf{A}_n^{(t)} - \eta_n^{(t)} \circledast \mathbf{G}_n^{(t)} \right).$$

- ▶ This is the adagrad scheme [Duchi et al., 2011]; see [Kolda and Hong, 2020] for using adam [Kingma and Ba, 2014] for entry sampling.

# Convergence



- ▶ **Setting:**  $l_1 = l_2 = l_3 = 100$ ,  $R = 10$ , NN latent factors.
- ▶ **Baselines:** AO-ADMM [Huang et al., 2016], APG [Xu and Yin, 2013].
- ▶ **Proposed:** BrasCPD (fine-tuned diminishing step size).  
AdaCPD (adaptive step size).
  - ▶  $|\mathcal{B}^{(t)}| = 9$  fibers sampled per iteration.

## Convergence Results - Details in [Fu et al., 2019b]

**Proposition 1.** Consider the case where  $h_n(\cdot) = 0$  for all  $n$ , that  $\alpha^{(t)}$  satisfies the Robinson-Monroe rule, and that the solution sequence is not unbounded. BrasCPD satisfies:

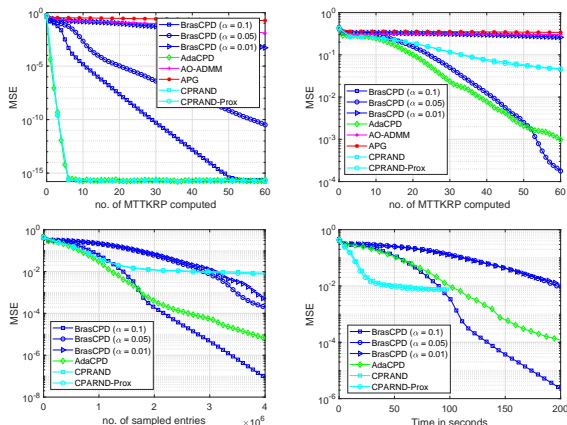
$$\liminf_{t \rightarrow \infty} \mathbb{E}[\|\nabla f(\boldsymbol{\theta}^{(t)})\|^2] = 0.$$

**Proposition 2.** In addition to the assumptions in Prop. 1, also assume that  $\Pr(\xi^{(t)} = n) = 1/N$  for all  $t$  and  $n$ . AdaCPD satisfies

$$\Pr\left(\liminf_{t \rightarrow \infty} \|\nabla f(\boldsymbol{\theta}^{(t)})\|^2 = 0\right) = 1.$$

**Proposition 3.** In addition to the assumptions in Prop. 1, also assume that the gradient estimation's variance diminishes when  $t \rightarrow \infty$ . Every limit point of BrasCPD's solution sequence is a stationary point.

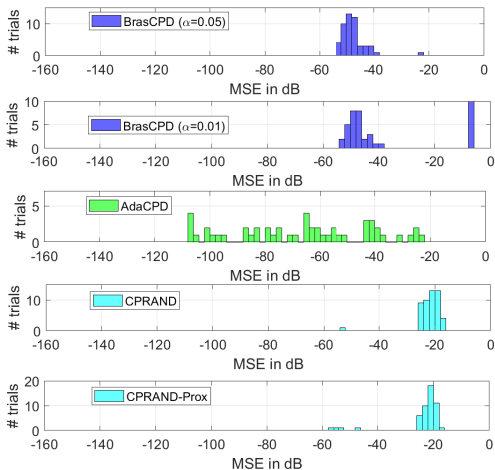
# More Results



**Setting:**  $I = 300, N = 3$ . Left upper:  $R = 10$ ; Others:  $R = 200$ .  
NN constraints. 50 trials.

- CPRAND [Battaglini et al., 2018]: fiber sampling and sketched LS (no constraint).

# More Results



**Setting:**  $I = 300, N = 3. R = 100.$  50 trials. NN constraints.

# Summary

- ▶ Stochastic optimization for tensor decomposition has become more important.
- ▶ Stochastic optimization's key considerations:
  - ▶ sampling schemes
  - ▶ step size scheduling (automatic/adaptive steps size is preferable)
  - ▶ constraints/reg.
  - ▶ convergence supports
- ▶ The block-randomized fiber sampling strategy seems to offer a promising solution package.

# Generalization for Non-Euclidean Losses

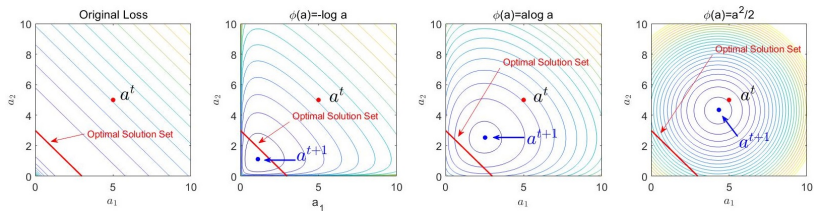
General problem of interest:

$$\min_{\text{model param.}} \text{dist}(\text{data}||\text{model}) + \left( \begin{array}{c} \text{penalty for} \\ \text{structure violation} \end{array} \right)$$

under structural constraints,

- ▶ How about  $\beta$ -divergence, KL-divergence, IS-divergence, Logistic loss, etc? [Cichocki et al., 2015, Chi and Kolda, 2012, Févotte et al., 2009]
  - ▶ used in integer, binary, and scaling-sensitive data analysis.
  - ▶ Entry sampling-based non-Euclidean CPD [Kolda and Hong, 2020, Hong et al., 2020].
    - ▶ SGD + Adam
- ▶ **Proposed:** Fiber sampling + stochastic **mirror descent**.

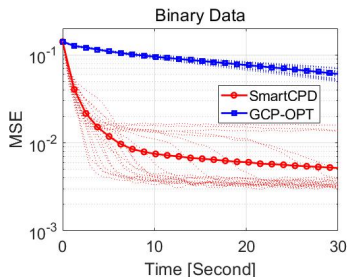
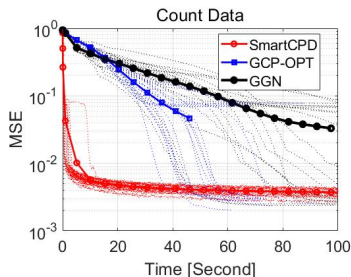
# More Results



MD adapts to loss function and constraint's geometry.



# SMD for CPD (SmartCPD)



- **Proposed:** SmartCPD - fiber sampling, block randomization, SMD
- **Baseline:** GCP-OPT [Hong et al., 2020, Kolda and Hong, 2020].  
GGN [Vandecappelle et al., 2020].
- **Freshly baked manuscript:**
  - W. Pu, S. Ibrahim, X. Fu, and M. Hong, “Stochastic Mirror Descent for Low-Rank Tensor Decomposition Under Non-Euclidean Losses”, arXiv:2104.14562 [stat.ML]

-  Battaglini, C., Ballard, G., and Kolda, T. G. (2018).  
A practical randomized CP tensor decomposition.  
39(2):876–901.
-  Beutel, A., Talukdar, P. P., Kumar, A., Faloutsos, C.,  
Papalexakis, E. E., and Xing, E. P. (2014).  
Flexifact: Scalable flexible factorization of coupled tensors on  
Hadoop.  
In *Proc. SIAM SDM 2014*, pages 109–117. SIAM.
-  Chi, E. C. and Kolda, T. G. (2012).  
On tensors, sparsity, and nonnegative factorizations.  
33(4):1272–1299.
-  Cichocki, A., Mandic, D., De Lathauwer, L., Zhou, G., Zhao,  
Q., Caiafa, C., and Phan, H.-A. (2015).  
Tensor decompositions for signal processing applications: From  
two-way to multiway component analysis.  
32(2):145–163.
-  Duchi, J., Hazan, E., and Singer, Y. (2011).

Adaptive subgradient methods for online learning and stochastic optimization.

[12\(Jul\):2121–2159.](#)



Févotte, C., Bertin, N., and Durrieu, J.-L. (2009).

Nonnegative matrix factorization with the itakura-saito divergence: With application to music analysis.

[Neural computation, 21\(3\):793–830.](#)



Fu, X., Huang, K., Sidiropoulos, N. D., and Ma, W.-K. (2019a).

Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications.

[IEEE Signal Process. Mag., 36\(2\):59–80.](#)



Fu, X., Ibrahim, S., Wai, H.-T., Gao, C., and Huang, K. (2019b).

Block-randomized stochastic proximal gradient for low-rank tensor factorization.

[arXiv preprint arXiv:1901.05529.](#)



Hong, D., Kolda, T. G., and Duersch, J. A. (2020).

Generalized canonical polyadic tensor decomposition.

*SIAM Review*, 62(1):133–163.



Huang, K., Sidiropoulos, N. D., and Liavas, A. P. (2016).

A flexible and efficient algorithmic framework for constrained matrix and tensor factorization.

*IEEE Trans. Signal Process.*, 64(19):5052–5065.



Kargas, N., Sidiropoulos, N. D., and Fu, X. (2018).

Tensors, learning, and 'Kolmogorov extension' for finite-alphabet random vectors.

*IEEE Trans. Signal Process.*, 66(18):4854–4868.



Kingma, D. P. and Ba, J. (2014).

Adam: A method for stochastic optimization.

*arXiv preprint arXiv:1412.6980*.



Kolda, T. G. and Bader, B. W. (2006).

Matlab tensor toolbox.

Technical report, Sandia National Laboratories.



Kolda, T. G. and Hong, D. (2020).

Stochastic gradients for large-scale tensor decomposition.

*SIAM Journal on Mathematics of Data Science*,  
2(4):1066–1095.



Liavas, A. P., Kostoulas, G., Lourakis, G., Huang, K., and Sidiropoulos, N. D. (2017).

Nesterov-based alternating optimization for nonnegative tensor factorization: Algorithm and parallel implementation.

*IEEE Trans. Signal Process.*, 66(4):944–953.



Lim, L.-H. and Comon, P. (2009).

Nonnegative approximations of nonnegative tensors.

23(7-8):432–441.



Lin, C.-J. (2007).

Projected gradient methods for nonnegative matrix factorization.

19(10):2756–2779.



Phipps, E. and Kolda, T. G. (2019).

Software for sparse tensor decomposition on emerging computing architectures.



Ravindran, N., Sidiropoulos, N. D., Smith, S., and Karypis, G. (2014).

Memory-efficient parallel computation of tensor and matrix products for big tensor decomposition.

*In 2014 48th Asilomar Conference on Signals, Systems and Computers*, pages 581–585. IEEE.



Razaviyayn, M., Hong, M., and Luo, Z.-Q. (2013).

A unified convergence analysis of block successive minimization methods for nonsmooth optimization.

23(2):1126–1153.



Sidiropoulos, N. D., De Lathauwer, L., Fu, X., Huang, K., Papalexakis, E. E., and Faloutsos, C. (2017).

Tensor decomposition for signal processing and machine learning.

*IEEE Trans. Signal Process.*, 65(13):3551–3582.



Smith, S., Ravindran, N., Sidiropoulos, N. D., and Karypis, G. (2015).

Splatt: Efficient and parallel sparse tensor-matrix multiplication.

*In 2015 IEEE International Parallel and Distributed Processing Symposium, pages 61–70.*



Tan, C., Ma, S., Dai, Y.-H., and Qian, Y. (2016).

Barzilai-borwein step size for stochastic gradient descent.  
*arXiv preprint arXiv:1605.04131.*



Vandecappelle, M., Vervliet, N., and De Lathauwer, L. (2020).

A second-order method for fitting the canonical polyadic decomposition with non-least-squares cost.  
*IEEE Transactions on Signal Processing, 68:4454–4465.*



Vervliet, N. and De Lathauwer, L. (2016).

A randomized block sampling approach to canonical polyadic decomposition of large-scale tensors.  
*IEEE J. Sel. Topics Signal Process., 10(2):284–295.*



Xu, Y. and Yin, W. (2013).

A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion.

6(3):1758–1789.



Yang, Y., Feng, Y., and Suykens, J. A. (2015).

A rank-one tensor updating algorithm for tensor completion.

*IEEE Signal Processing Letters*, 22(10):1633–1637.



Yeredor, A. and Haardt, M. (2019).

Maximum likelihood estimation of a low-rank probability mass tensor from partial observations.

*IEEE Signal Process. Lett.*, 26(10):1551–1555.



Zeiler, M. D. (2012).

Adadelta: an adaptive learning rate method.

*arXiv preprint arXiv:1212.5701*.