

# Non-negative Tensor Decompositions for Unsupervised Learning

Joint work w/ Naomi Graham<sup>1</sup>, Joel Saylor<sup>2</sup>, & Michael Friedlander<sup>1,3</sup> 2023 SIAM PNW Conference at Western Washington University, Bellingham, WA

Nicholas Richardson<sup>3</sup>

October 14, 2023

<sup>1</sup>UBC Computer Science <sup>2</sup>UBC Earth, Ocean, and Atmospheric Sciences <sup>3</sup>UBC Mathematics



THE UNIVERSITY OF BRITISH COLUMBIA



## Table of Contents

► Overview

Tensor Models

Computing Tensor Decompositions

Special Considerations

Conclusion



## The Setting

- Do you have a lot of data?
- A function that depends on many independent variables?
- High-dimensional samples?
- Want to compress the data?
- Or find relationships between samples/variables?





### Benefits Overview

- Less storage than original data
- Interpretable results
- There exist easy to implement algorithms with convergence guarantees
- Unlike supervised learning, no "training" required
- None or only a few hyperparameters to tune



Table of Contents Tensor Models

► Overview

► Tensor Models

Computing Tensor Decompositions

Special Considerations

► Conclusion



#### The Data Tensor Tensor Models



Figure: Order-3 Tensor

<sup>1</sup>Kolda and Bader 2009

• Order-N tensor<sup>1</sup>

 $Y \in \mathbb{R}^{I_1 imes \cdots imes I_N}$ 

- Often consider additional constraints
  - non-negative:  $Y_{i_1...i_N} \in \mathbb{R}_+$
  - probability values:  $0 \leq Y_{i_1...i_N} \leq 1$
  - distributions:  $1 = \sum_{j} Y_{i_1...i_N}$
  - binary:  $Y_{i_1...i_N} \in \{0, 1\}$



### Tucker Decomposition<sup>1</sup> Tensor Models

• Factorize Y into a core tensor  $G \in \mathbb{R}^{R_1 imes \cdots imes R_N}$  and matrices  $A^n \in \mathbb{R}^{I_n imes R_n}$ 

$$Y = G imes_1 A^1 imes_2 \cdots imes_N A^N$$
 $Y_{i_1...i_N} = \sum_{r_1...r_N} G_{r_1...r_N} A^1_{i_1r_1} \cdots A^N_{i_Nr_N}$ 

- e.g. compress  $1000^3$  (4 GB) to  $100^3 + 3 \cdot 100 \cdot 1000$  (5.2 MB)



<sup>1</sup>Tucker 1966



### Special Cases of Tucker Decomposition Tensor Models

- Canonical polyadic decomposition<sup>1</sup> (CP): G = I and  $R_1, \ldots, R_N = R$
- Smallest *R* is the tensor rank

$$egin{aligned} Y &= I imes_1 A^1 imes_2 \cdots imes_N A^N \ Y_{i_1 \dots i_N} &= \sum_r A^1_{i_1 r} \cdots A^N_{i_N r} \ Y_{i_N r} \ Y_{i_$$

e.g. low-rank decomp:  $Y = AB^{\top}$ 

- Tucker- $n^1$ :  $A^{n+1}$ , ...,  $A^N = I$  (possibly different sizes!)
  - $Y = G \times_1 A^1 \times_2 \cdots \times_n A^n$ e.g. Tucker-1/matrix-tensor factor<sup>2</sup>:  $Y = G \times_1 A = AG$

<sup>1</sup>Kolda and Bader 2009 <sup>2</sup>Richardson, Graham, et al. n.d.



### CP Decomposition Example 1 Tensor Models

- Decompose one bass line (STFT) into Fourier Transform & Amplitude Envelope<sup>1</sup>
  - Order N = 2, rank R = 1



<sup>1</sup>Richardson n.d.



### CP Decomposition Example 2

#### Tensor Models

- Decompose multiple bass lines into notes & envelopes for each song<sup>1</sup>
  - Order N = 3, rank R = 11



<sup>1</sup>Richardson n.d.



## Tucker-1 Example

- Decompose multiple sink's feature densities into linear combination of latent source densities<sup>1</sup>
  - Order 3, rank 3



<sup>&</sup>lt;sup>1</sup>Richardson, Graham, et al. n.d.



### Table of Contents Computing Tensor Decompositions

Overview

Tensor Models

Computing Tensor Decompositions

Special Considerations

► Conclusion



### The Optimization Problem Computing Tensor Decompositions

• May try to solve the model exactly

$$Y = f(A^1, \dots, A^N)$$
 (e.g. NMF:  $Y = WH$ )

• With noise, or imperfect modelling, better to find "closest" fit

$$\min_{A^n\in\mathcal{C}^n}D(Y,\hat{Y})$$
 s.t.  $\hat{Y}=f(A^1,\ldots,A^N)$ 

• Define objective function  $D(\mathbf{Y}, \hat{\mathbf{Y}}) = F(A^1, \dots, A^N)$ 

$$\begin{array}{l} - \hspace{0.1 cm} \text{e.g. Least Squares: } \|Y - \hat{Y}\|_{F}^{2} \\ - \hspace{0.1 cm} \text{e.g. KL-Divergence: } \sum_{i_{1}...i_{N}} Y_{i_{1}...i_{N}} \log \left( \frac{Y_{i_{1}...i_{N}}}{\hat{Y}_{i_{1}...i_{N}}} \right) \end{array}$$



### **Problems With Naive Methods**

**Computing Tensor Decompositions** 

• (Projected) Full Gradient Descent:

$$\mathcal{A} \leftarrow \arg \min_{\mathcal{B} \in \mathcal{C}} \langle 
abla F(\mathcal{A}), \mathcal{B} - \mathcal{A} 
angle + rac{1}{2lpha} \|\mathcal{B} - \mathcal{A}\|_F^2$$
  
 $(\mathcal{A}^1, \dots, \mathcal{A}^N) \leftarrow P_{\mathcal{C}} \left( (\mathcal{A}^1, \dots, \mathcal{A}^N) - lpha 
abla F(\mathcal{A}^1, \dots, \mathcal{A}^N) 
angle 
ight)$ 

- Non-expensive updates, but F not convex,  $\nabla F$  not Lipschitz
- Alternating Least-Squares:

$$A^n \leftarrow \arg \min_{A \in \mathcal{C}} ||Y - f(A^1, \dots, A^{n-1}, A, A^{n+1}, \dots, A^N)||_F^2$$

- f linear in  $A^n$ , but waste time fully optimizing  $A^n$  every step



### **Block Coordinate Descent (BCD) Algorithm**<sup>2</sup>

Computing Tensor Decompositions

• Combine into BCD<sup>1</sup>:

$$\mathbf{A}^n \leftarrow rg \min_{\mathbf{A} \in \mathcal{C}} \langle 
abla_{\mathbf{A}^n} F, \mathbf{A} - \mathbf{A}^n 
angle + rac{1}{2lpha} \|\mathbf{A} - \mathbf{A}^n\|_F^2$$

$$\mathbf{A}^{n} \leftarrow P_{\mathcal{C}} \left( \mathbf{A}^{n} - \alpha \nabla_{\mathbf{A}^{n}} F \right)$$

- Alternately update the factors, but don't fully optimize every iteration
- Can choose  $\alpha = 1/L_n$ 
  - *L* is the Lipshitz constant of  $\nabla F(A^1, \dots, A^{n-1}, \cdot, A^{n+1}, \dots, A^N)$



Figure: Example minimization with BCD (blue), Gradient Descent (green), Alternating Least Squares (orange)

<sup>&</sup>lt;sup>1</sup>Xu and Yin 2013

<sup>11 &</sup>lt;sup>2</sup>https://github.com/njericha/Sediment-Source-Analysis.jl



Table of Contents Special Considerations

Overview

Tensor Models

Computing Tensor Decompositions

► Special Considerations

Conclusion



### **Convergence** Special Considerations

• BCD converges to (global) Nash point  $(A^1, \ldots, A^N)$ :

$$A^n = rg \min_{A \in \mathcal{C}^n} F(A^1, \dots, A^{n-1}, A, A^{n+1}, \dots, A^N), \quad n = 1, \dots, N$$

• "Cannot improve objective by updating one block"

### Nash point condition<sup>2</sup>

Let  $F(\cdot,B),F(A,\cdot)$  be convex, and  $(A_0,B_0)\in\mathcal{C}=\mathcal{C}_A imes\mathcal{C}_B.$  Then,

$$\mathbf{0} \in \partial(F + \delta_{\mathcal{C}})(A_0, B_0) \iff \begin{array}{c} F(A_0, B_0) \leq F(A, B_0) \; \forall A \in \mathcal{C}_A \\ F(A_0, B_0) \leq F(A_0, B) \; \forall B \in \mathcal{C}_B \end{array}$$

• When *F* is differentiable and *C* is convex:

$$\mathbf{0} \in \partial(F + \delta_{\mathcal{C}})(A, B) \iff -\nabla F \in N_{\mathcal{C}}(A, B)$$

<sup>12</sup> <sup>1</sup>Xu and Yin 2013 <sup>2</sup>Richardson, Graham, et al. n.d.



### Selecting the Rank Special Considerations

- Option 1: Use information about your physical system
  - e.g. decompose piano audio into notes and amplitudes
  - -R = 88 since there are 88 keys
- Option 2: compute the "best bank for your buck"<sup>1</sup>
  - Solve the model for all ranks  $r=1,\ldots,I$
  - Compute final objective value F(R)
  - Select point of maximum curvature

$$R = rg \max_{r} \kappa(r) := rac{F''(r)}{(1+F'(r)^2)^{3/2}}$$



Figure: Typical final error F(r) vs rank r plot

<sup>1</sup>Satopaa et al. 2011



## Uniqueness & Scaling

- These models are not (usually) unique
- Ex.  $Y = AB^{\top} = (AC)(C^{-1}B^{\top})$  for invertable C
- Fix a scaling on factors: set  $\|A^n\| = c_n$  or  $\sum_j A_{..j..} = c_n$  for...

... all but one factor

 $\ldots$  all factors & add scaling parameter  $\lambda$ 

e.g.: 
$$Y_{ij} = \sum_{r} A_{ir} B_{jr}$$
 s.t.  $||B_{:r}||_2 = 1$  e.g.:  $Y_{ij} = \sum_{r} \lambda_r A_{ir} B_{jr}$  s.t.  $||A_{:r}||_2, ||B_{:r}||_2 = 1$ 

- Enforce through:
  - constraint (projection)
  - rescale at the end/each iteration
- Still only unique up to permutations of rows/columns/fibres



## Table of Contents

Overview

Tensor Models

Computing Tensor Decompositions

Special Considerations

### ► Conclusion



#### Summary Conclusion

- Looked at various tensor decomposition models
- Optimization methods to solve them
- Practical considerations



# Future Directions for Additional Compression

• Factorize the core symmetrically<sup>1</sup> (Extend Tucker)

$$Y_{ijk} = \sum_{r_1, r_2, r_3} \sum_{\substack{p,q,s=1\\p,q,s=1}}^{R} B_{r_1qs}^1 B_{pr_2s}^2 B_{pqr_3}^3 A_{ir_1}^1 A_{jr_2}^2 A_{kr_3}^3$$
  
• Tensor Trains<sup>2</sup> (Extend Tucker-2)  

$$Y_{i_1...i_N} = \sum_{r_1, r_N} A_{i_1r_1}^1 \sum_{\substack{j_2...j_{N-2}\\p_2...j_{N-2}}} A_{r_1i_2j_2}^2 A_{j_2i_3j_3}^3 \dots A_{j_{N-2}i_{N-1}r_2}^{n-1} A_{i_Nr_N}^N$$

• Factorize the matrices (Extend CP)

$$Y_{ij} = \sum_{r} A_{ir} \underbrace{\sum_{k} T_{jrk} b_{k}}_{B_{ir}}$$

<sup>1</sup>Qi et al. 2020 <sup>2</sup>Oseledets 2011



### References

- Kolda, Tamara G. and Brett W. Bader (Aug. 2009). "Tensor Decompositions and Applications". In: *SIAM Review* 51.3, pp. 455–500.
- Luo, Yuan, Fei Wang, and Peter Szolovits (May 2017). "Tensor Factorization toward Precision Medicine". In: *Briefings in Bioinformatics* 18.3, pp. 511–514.
- Oseledets, I. V. (Jan. 2011). "Tensor-Train Decomposition". In: SIAM Journal on Scientific Computing 33.5, pp. 2295–2317.



- Qi, Liqun et al. (Mar. 2020). Triple Decomposition and Tensor Recovery of Third Order Tensors.
- Richardson, Nicholas (n.d.). A Consistant Framework for Non-negative Tensor Models and Algorithms with Applications.
- Richardson, Nicholas, Naomi Graham, et al. (n.d.). Non-Negative Matrix-Tensor Factorization for Sediment Source Analysis.

Richardson, Nicholas, Hayden Schaeffer, and Giang Tran

<sup>17</sup> (June 2023). "SRMD: Sparse Random Mode Decomposi-

tion". In: Communications on Applied Mathematics and Computation.

- Satopaa, Ville et al. (June 2011). "Finding a "Kneedle" in a Haystack: Detecting Knee Points in System Behavior". In: 2011 31st International Conference on Distributed Computing Systems Workshops, pp. 166–171.
- Sundell, Kurt E. et al. (2022). "Crustal Thickening of the Northern Central Andean Plateau Inferred From Trace Elements in Zircon". In: *Geophysical Research Letters* 49.3, e2021GL096443.
- Suvanjanprasai (n.d.). Sample Images from MNIST Test Dataset.
- Tucker, Ledyard R. (Sept. 1966). "Some Mathematical Notes on Three-Mode Factor Analysis". In: *Psychometrika* 31.3, pp. 279–311.
- Xu, Yangyang and Wotao Yin (Jan. 2013). "A Block Coordinate Descent Method for Regularized Multiconvex Optimization with Applications to Nonnegative Tensor Factorization and Completion". In: SIAM Journal on Imaging Sciences 6.3, pp. 1758–1789.



Thank you for listening! Any questions?