# Neural Nonnegative Matrix Factorization for Hierarchical Multilayer Topic Modeling

Jamie Haddock CAMSAP 2019, December 16, 2019

Computational and Applied Mathematics UCLA



joint with Mengdi Gao, Denali Molitor, Deanna Needell, Eli Sadovnik, Tyler Will, Runyu Zhang











k: user chosen parameter



k: user chosen parameter

 $\triangleright$  nonconvex in **A** and **S**, NP-hard [Vavasis '08]



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}} \|\mathbf{X} - \mathbf{AS}\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{X} &\in \mathbb{R}_{\geq 0}^{N \times M} \text{: data matrix} \\ \mathbf{A} &\in \mathbb{R}_{\geq 0}^{N \times k} \text{: features matrix} \\ \mathbf{S} &\in \mathbb{R}_{\geq 0}^{k \times M} \text{: coefficients matrix} \\ k \text{: user chosen parameter} \end{split}$$

#### **Problem Challenges:**

- ▷ nonconvex in A and S, NP-hard [Vavasis '08]
- ▷ interpretability of factors dependent upon k



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}} \|\mathbf{X} - \mathbf{AS}\|_{F}^{2}$$

#### **Problem Setup:**

 $\begin{array}{l} \mathbf{X} \in \mathbb{R}_{\geq 0}^{N \times M}: \text{ data matrix} \\ \mathbf{A} \in \mathbb{R}_{\geq 0}^{N \times k}: \text{ features matrix} \\ \mathbf{S} \in \mathbb{R}_{\geq 0}^{k \times M}: \text{ coefficients matrix} \\ k: \text{ user chosen parameter} \end{array}$ 

#### **Problem Challenges:**

- ▷ nonconvex in A and S, NP-hard [Vavasis '08]
- interpretability of factors dependent upon k

### Methods:

▷ low-rank approximation

- ▷ low-rank approximation
- ▷ clustering

- ▷ low-rank approximation
- ▷ clustering
- ▷ topic modeling

- ▷ low-rank approximation
- ▷ clustering
- ▷ topic modeling
- ▷ feature extraction

- ▷ low-rank approximation
- ▷ clustering
- ▷ topic modeling
- ▷ feature extraction

### Methods:

> multiplicative updates

- ▷ low-rank approximation
- ▷ clustering
- ▷ topic modeling
- ▷ feature extraction

- > multiplicative updates
- alternating nonnegative least squares

- ▷ low-rank approximation
- ▷ clustering
- ▷ topic modeling
- ▷ feature extraction

- > multiplicative updates
- alternating nonnegative least squares
- $\triangleright$  many others

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

**Problem Setup:** 

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

Problem Setup:  $\mathbf{Y} \in \{0,1\}_{>0}^{P \times M}: \text{ label matrix}$ 

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

 $\mathbf{Y} \in \{0,1\}_{\geq 0}^{P imes M}$ : label matrix *P*: number of classes

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{Y} &\in \{0,1\}_{\geq 0}^{P \times M} \text{: label matrix} \\ P \text{: number of classes} \\ \mathbf{W} &\in \{0,1\}_{\geq 0}^{N \times M} \text{: data indicator} \end{split}$$

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{Y} &\in \{0,1\}_{\geq 0}^{P \times M} \text{: label matrix} \\ P \text{: number of classes} \\ \mathbf{W} &\in \{0,1\}_{\geq 0}^{N \times M} \text{: data indicator} \\ \mathbf{L} &\in \{0,1\}_{\geq 0}^{P \times M} \text{: label indicator} \end{split}$$

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{Y} &\in \{0,1\}_{\geq 0}^{P \times M}: \text{ label matrix } \\ P: \text{ number of classes } \\ \mathbf{W} &\in \{0,1\}_{\geq 0}^{N \times M}: \text{ data indicator } \\ \mathbf{L} &\in \{0,1\}_{\geq 0}^{P \times M}: \text{ label indicator } \\ \lambda: \text{ user defined hyperparameter } \end{split}$$

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{Y} &\in \{0,1\}_{\geq 0}^{P \times M} \text{: label matrix} \\ P \text{: number of classes} \\ \mathbf{W} &\in \{0,1\}_{\geq 0}^{N \times M} \text{: data indicator} \\ \mathbf{L} &\in \{0,1\}_{\geq 0}^{P \times M} \text{: label indicator} \\ \lambda \text{: user defined hyperparameter} \end{split}$$

#### **Problem Advantages:**

▷ use of label information

Goal: Incorporate known label information into problem.



$$\min_{\mathbf{A} \in \mathbb{R}^{N \times k}_{\geq 0}, \mathbf{S} \in \mathbb{R}^{k \times M}_{\geq 0}, \mathbf{B} \in \mathbb{R}^{P \times k}_{\geq 0}} \|\mathbf{W} \odot (\mathbf{X} - \mathbf{AS})\|_{F}^{2} + \lambda \|\mathbf{L} \odot (\mathbf{Y} - \mathbf{BS})\|_{F}^{2}$$

#### **Problem Setup:**

$$\begin{split} \mathbf{Y} &\in \{0,1\}_{\geq 0}^{P \times M}: \text{ label matrix} \\ P: \text{ number of classes} \\ \mathbf{W} &\in \{0,1\}_{\geq 0}^{N \times M}: \text{ data indicator} \\ \mathbf{L} &\in \{0,1\}_{\geq 0}^{P \times M}: \text{ label indicator} \\ \lambda: \text{ user defined hyperparameter} \end{split}$$

- ▷ use of label information
- can extend multiplicative updates method to SSNMF

**Goal:** Discover hierarchical topic structure within **X**.

**Problem Setup:** 

**Problem Challenges:** 

Goal: Discover hierarchical topic structure within X.



 $\boldsymbol{\mathsf{X}} \approx \boldsymbol{\mathsf{A}}^{(0)} \boldsymbol{\mathsf{A}}^{(1)} \dots \boldsymbol{\mathsf{A}}^{(\mathcal{L})} \boldsymbol{\mathsf{S}}^{(\mathcal{L})}$ 

Goal: Discover hierarchical topic structure within X.



$$\begin{split} \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{S}^{(0)} \\ \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \mathbf{S}^{(1)} \end{split}$$

Problem Setup:

 $\triangleright k^{(0)}, k^{(1)}, \dots, k^{(\mathcal{L})}$ : user defined parameters

**Problem Challenges:** 

 $\boldsymbol{\mathsf{X}} \approx \boldsymbol{\mathsf{A}}^{(0)} \boldsymbol{\mathsf{A}}^{(1)} \dots \boldsymbol{\mathsf{A}}^{(\mathcal{L})} \boldsymbol{\mathsf{S}}^{(\mathcal{L})}$ 

Goal: Discover hierarchical topic structure within X.



$$\begin{split} \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{S}^{(0)} \\ \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \mathbf{S}^{(1)} \end{split}$$

 $\mathbf{X} \approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \qquad \mathbf{A}^{(\mathcal{L})} \mathbf{S}^{(\mathcal{L})}$ 

**Problem Setup:** 

▷ k<sup>(0)</sup>, k<sup>(1)</sup>,..., k<sup>(L)</sup>: user defined parameters
▷ k<sup>(ℓ)</sup>: supertopics collecting k<sup>(ℓ-1)</sup> subtopics
Problem Challenges:

5

**Goal:** Discover hierarchical topic structure within **X**.



 $\mathbf{X} \approx \mathbf{A}^{(0)} \mathbf{S}^{(0)}$  $\mathbf{X} \approx \mathbf{\Delta}^{(0)} \mathbf{\Delta}^{(1)} \mathbf{S}^{(1)}$  **Problem Setup:** 

 $\triangleright k^{(0)}, k^{(1)}, \ldots, k^{(\mathcal{L})}$ : user defined parameters  $\triangleright k^{(\ell)}$ : supertopics collecting  $k^{(\ell-1)}$  subtopics **Problem Challenges:**  $\mathbf{X} \approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \dots \mathbf{A}^{(\mathcal{L})} \mathbf{S}^{(\mathcal{L})} \qquad \triangleright \{k^{(i)}\} \text{ must be chosen}$ 

Goal: Discover hierarchical topic structure within X.



$$\begin{split} \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{S}^{(0)} \\ \mathbf{X} &\approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \mathbf{S}^{(1)} \end{split}$$

 $\mathbf{X} \approx \mathbf{A}^{(0)} \mathbf{A}^{(1)} \dots \mathbf{A}^{(\mathcal{L})} \mathbf{S}^{(\mathcal{L})}$ 

#### **Problem Setup:**

▷  $k^{(0)}, k^{(1)}, \dots, k^{(\mathcal{L})}$ : user defined parameters ▷  $k^{(\ell)}$ : supertopics collecting  $k^{(\ell-1)}$  subtopics

### Problem Challenges:

▷  $\{k^{(i)}\}$  must be chosen

▷ error propagates through layers



- $\triangleright$  [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate

- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF

- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF
- ▷ [Le Roux, Hershey, Weninger '15]
  - introduces NMF backpropagation algorithm with "unfolding" (no hierarchy)

- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF
- ▷ [Le Roux, Hershey, Weninger '15]
  - introduces NMF backpropagation algorithm with "unfolding" (no hierarchy)
- ▷ [Sun, Nasrabadi, Tran '17]
  - similar method lacking nonnegativity constraints

Goal: Develop true backpropagation algorithm for hNMF model.

Goal: Develop true backpropagation algorithm for hNMF model.

Regard the A matrices as independent variables, determine the S matrices from the A matrices.

Goal: Develop true backpropagation algorithm for hNMF model.

- Regard the A matrices as independent variables, determine the S matrices from the A matrices.
- $\triangleright \text{ Define } q(\mathbf{X}, \mathbf{A}) := \operatorname{argmin}_{S \ge 0} \|\mathbf{X} \mathbf{AS}\|_{F}^{2}.$

Goal: Develop true backpropagation algorithm for hNMF model.

- Regard the A matrices as independent variables, determine the S matrices from the A matrices.
- $\triangleright \text{ Define } q(\mathbf{X}, \mathbf{A}) := \operatorname{argmin}_{S \ge 0} \|\mathbf{X} \mathbf{AS}\|_{F}^{2}.$



Goal: Develop true backpropagation algorithm for hNMF model.

- Regard the A matrices as independent variables, determine the S matrices from the A matrices.
- $\triangleright \text{ Define } q(\mathbf{X}, \mathbf{A}) := \operatorname{argmin}_{S \ge 0} \|\mathbf{X} \mathbf{AS}\|_{F}^{2}.$



▷ Pin the values of **S** to those of **A** by recursively setting  $S^{(\ell)} := q(S^{(\ell-1)}, A^{(\ell)}).$ 

Goal: Develop true backpropagation algorithm for hNMF model.

- Regard the A matrices as independent variables, determine the S matrices from the A matrices.
- $\triangleright \text{ Define } q(\mathbf{X}, \mathbf{A}) := \operatorname{argmin}_{S \ge 0} \|\mathbf{X} \mathbf{AS}\|_{F}^{2}.$



- ▷ Pin the values of **S** to those of **A** by recursively setting  $\mathbf{S}^{(\ell)} := q(\mathbf{S}^{(\ell-1)}, \mathbf{A}^{(\ell)}).$
- ▷ Can we compute derivatives and backpropagate?

# **Neural NMF Backpropagation**





▷ Differentiate q function and apply chain rule.



- ▷ Differentiate q function and apply chain rule.
- Flexible to cost function (e.g., supervision).



- ▷ Differentiate q function and apply chain rule.
- Flexible to cost function (e.g., supervision).
- Backpropagate and update all A matrices simultaneously via GD or SGD.

#### Method 1 Neural NMF

**Require:** data matrix  $\mathbf{X} \in \mathbb{R}^{N \times M}$ , number of layers  $\mathcal{L}$ , step size  $\gamma$ , cost function C, initial matrices  $\mathbf{A}^{(i)}$  for  $i = 0, ..., \mathcal{L}$  **procedure** FORWARDPROPAGATION $(\mathbf{A}^{(0)}, ..., \mathbf{A}^{(\mathcal{L})})$  **for**  $i := 0...\mathcal{L}$  **do**  $\mathbf{S}^{(i)} \leftarrow q(\mathbf{A}^{(i)}, \mathbf{S}^{(i-1)})$ 

ForwardPropagation( $\mathbf{A}^{(0)}, \dots, \mathbf{A}^{(\mathcal{L})}$ ) while not converged do for  $i := 0...\mathcal{L}$  do  $\mathbf{A}^{(i)} \leftarrow \mathbf{A}^{(i)} - \gamma * \frac{\partial C}{\partial \mathbf{A}^{(i)}}$   $\triangleright$  Gradient descent  $\mathbf{A}^{(i)} \leftarrow \mathbf{A}^{(i)}_{+}$   $\triangleright$  Project onto positive orthant

ForwardPropagation( $\mathbf{A}^{(0)}, \dots, \mathbf{A}^{(\mathcal{L})}$ )



 $\triangleright$  unsupervised reconstruction with two-layer structure  $(k^{(0)} = 9, k^{(1)} = 4)$ 



▷ semisupervised reconstruction (40% labels) with three-layer structure  $(k^{(0)} = 9, k^{(1)} = 4, k^{(2)} = 2)$ 

Note that despite reconstruction error increasing as layers increase (since the final rank decreases), the topic structure can be resolved from the intermediate factorizations.



▷ unsupervised reconstruction with two-layer structure  $(k^{(0)} = 9, k^{(1)} = 4)$ 

Note that despite reconstruction error increasing as layers increase (since the final rank decreases), the topic structure can be resolved from the intermediate factorizations.



▷ semisupervised reconstruction (40% labels) with three-layer structure  $(k^{(0)} = 9, k^{(1)} = 4, k^{(2)} = 2)$ 

#### Table 1: Reconstruction error / classification accuracy

	Layers	Hier. NMF	Deep NMF	Neural NMF
Unsuper.	1	0.053	0.031	0.029
	2	0.399	0.414	0.310
	3	0.860	0.838	0.492
Semisuper.	1	0.049 / 0.933	0.031 / 0.947	0.042 / 1
	2	0.374 / 0.926	0.394 / 0.911	0.305 / 1
	3	0.676 / 0.930	0.733 / 0.930	0.496 / 0.990
Supervised	1	0.052 / 0.960	0.042 / 0.962	0.042 / 1
	2	0.311 / 0.984	0.310 / 0.984	0.307 / 1
	3	0.495 / 1	0.494 / 1	0.498 / 1

presented a novel method for multilayer NMF that incorporates the backpropagation technique from deep learning to minimize error accumulation

- presented a novel method for multilayer NMF that incorporates the backpropagation technique from deep learning to minimize error accumulation
- exhibited preliminary tests on toy datasets showing the proposed method outperforms existing multilayer NMF algorithms.

- presented a novel method for multilayer NMF that incorporates the backpropagation technique from deep learning to minimize error accumulation
- exhibited preliminary tests on toy datasets showing the proposed method outperforms existing multilayer NMF algorithms.

▷ compare our method and others on various datasets to find precise regimes in which we offer improvement

- presented a novel method for multilayer NMF that incorporates the backpropagation technique from deep learning to minimize error accumulation
- exhibited preliminary tests on toy datasets showing the proposed method outperforms existing multilayer NMF algorithms.

- ▷ compare our method and others on various datasets to find precise regimes in which we offer improvement
- $\triangleright~$  extend to method for hierarchical nonnegative tensor factorization

# Questions?

- [1] Jennifer Flenner and Blake Hunter. A deep non-negative matrix factorization neural network, 2018. Unpublished.
- [2] Jonathan Le Roux, John R Hershey, and Felix Weninger. Deep nmf for speech separation. In Int. Conf. Acoust. Spee., pages 66–70. IEEE, 2015.
- [3] Daniel D. Lee and H. Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. <u>Nature</u>, 401:788–791, 1999.
- [4] H. Lee, J. Yoo, and S. Choi. Semi-supervised nonnegative matrix factorization. IEEE Signal Proc. Let., 17(1):4–7, Jan 2010.
- [5] Xiaoxia Sun, Nasser M. Nasrabadi, and Trac D. Tran. Supervised multilayer sparse coding networks for image classification. CoRR, abs/1701.08349, 2017.
- [6] George Trigeorgis, Konstantinos Bousmalis, Stefanos Zafeiriou, and Björn W Schuller. A deep matrix factorization method for learning attribute representations. IEEE T. Pattern Anal., 39(3):417–429, 2016.