Accelerating Non-negative Matrix Factorization Algorithms using Extrapolation, and more i.e. How to  $\min_{\mathbf{W},\mathbf{H}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F$  subject to  $\mathbf{W} \ge 0, \mathbf{H} \ge 0$ 

Andersen Ang

Mathématique et recherche opérationnelle, UMONS, Belgium Email: manshun.ang@umons.ac.be Homepage: angms.science

> PANAMA, {INRIA,IRISA}, Rennes, France Feburary 12, 2019



## Find $(\mathbf{W}, \mathbf{H})$ numerically

- Variations on BCD
- A-HALS
- Projected Gradient Update and the Multiplicative update

## 3) Find $(\mathbf{W},\mathbf{H})$ numerically fast : acceleration via extrapolation

- Recall : acceleration in single variable problem
- Accelerating NMF algorithms using extrapolation
- 4 Convergence of the algorithms
  - Application of PALM on NMF
  - Convergence condition of PALM

# Outline



- Find  $(\mathbf{W}, \mathbf{H})$  numerically
- Variations on BCD
- A-HALS
- Projected Gradient Update and the Multiplicative update

# 3) Find $(\mathbf{W},\mathbf{H})$ numerically fast : acceleration via extrapolation

- Recall : acceleration in single variable problem
- Accelerating NMF algorithms using extrapolation
- Convergence of the algorithms
  - Application of PALM on NMF
  - Convergence condition of PALM

# Non-negative Matrix Factorization (NMF)

Given :

- A matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ .
- A positive integer  $r \in \mathbb{N}$ .

Given :

- A matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ .
- A positive integer  $r \in \mathbb{N}$ .

Find :

• Matrices  $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$  such that  $\mathbf{X} = \mathbf{W}\mathbf{H}$ .

Given :

- A matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ .
- A positive integer  $r \in \mathbb{N}$ .

Find :

- Matrices  $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$  such that  $\mathbf{X} = \mathbf{W}\mathbf{H}$ .
- Important : everything is non-negative.



# Exact and approximate NMF

Given  $(\mathbf{X} \in \mathbb{R}^{m \times n}_+, r \in \mathbb{N})$ , find  $(\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+)$ s.t.  $\mathbf{X} = \mathbf{W}\mathbf{H}$  is called *exact NMF*.

## It is NP-hard (Vavasis, 2007).

Vavasis, "On the complexity of nonnegative matrix factorization", SIAM J. Optim.

# Exact and approximate NMF

Given  $(\mathbf{X} \in \mathbb{R}^{m \times n}_+, r \in \mathbb{N})$ , find  $(\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+)$ s.t.  $\mathbf{X} = \mathbf{W}\mathbf{H}$  is called *exact NMF*.

## It is NP-hard (Vavasis, 2007).

Vavasis, "On the complexity of nonnegative matrix factorization", SIAM J. Optim.

# This talk : (Low-rank) approximate NMF $\mathbf{X} \approx \mathbf{WH}, \ 1 \leq r \leq \min\{m, n\}.$



# Find $(\mathbf{W},\mathbf{H})$ numerically

Given  $(\mathbf{X} \in \mathbb{R}^{m \times n}_+, 1 \le r \le \min\{m, n\})$ , find  $(\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+)$  s.t.  $\mathbf{X} \approx \mathbf{W}\mathbf{H}$  via solving  $[\mathbf{W}, \mathbf{H}] = \underset{\mathbf{W} > \mathbf{0}, \mathbf{H} > \mathbf{0}}{\arg \min} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F.$ 

- $\bullet$  Minimizing the distance between  ${\bf X}$  and the approximator  ${\bf WH}$  in F-norm^†.
- $\geq$  is element-wise (not positive semi-definite).
- Such minimization problem is
  - Bi-variate : two variables
  - $\blacktriangleright$  Non-smooth : on the boundary between  ${\rm I\!R}_+$  and  ${\rm I\!R}_-$
  - Non-convex
  - III-posed and NP-hard (Vavasis, 2007)

†This talk does not consider other distance functions.

Solving the minimization problem

$$[\mathbf{W}, \ \mathbf{H}] = \underset{\mathbf{W} \ge \mathbf{0}, \mathbf{H} \ge \mathbf{0}}{\arg\min} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F},$$

**Keywords** : Numerical optimization, Continuous optimization, Algorithm, Convergence, Non-convex, Nesterov's Acceleration, Extrapolation

**Non-keywords** : Sparsity, Regularization, Applications of NMFs, Extended Formulations, Separability, Non-negative rank

# 4 slides on why NMF c'est bon bon

# For non-NMF people : why NMF ?

# Interpretability

NMF beats similar tools (PCA, SVD, ICA) due to the interpretability

on non-negative data.

# Model correctness

NMF can find ground truth (under certain conditions).

# • Mathematical curiosity

NMF is related to some serious problems in mathematics.

# • My boss tell me to do it.

# Why NMF - Hyper-spectral image application (1/2)

NMF gives good *unsupervised* image segmentation<sup>1</sup>



Figure: Hyper-spectral image decomposition. Figure from (Zhu,2014).

Zhu, F. et al., Spectral unmixing via data-guided sparsity. IEEE Trans. Image Processing

# Comment est-ce possible ?!

<sup>&</sup>lt;sup>1</sup>Modern fancy name : "super resolution"

# Why NMF - Hyper-spectral image application (2/2)



Figure: Hyper-spectral imaging. Figure modified from the slide of Nicolas Gillis. 14/119

# Why NMF - other examples

#### Application side

- Spectral unmixing in analytical chemistry (one of the earliest work)
- Representation learning on human face (the work that popularizes NMF)
- Topic modeling in text mining
- Probability distribution application on identification of Hidden Markov Model
- Bioinformatics : gene expression
- Time-frequency matrix decompositions for neuroinformatics
- (Non-negative) Blind source separation
- (Non-negative) Data compression
- Speech denoising
- Recommender system
- Face recognition
- Video summarization
- Radio
- Forensics
- Art work conservation (identify true color used in painting)
- Medical imaging image processing on small object
- Mid-infrared astronomy image processing on large object
- Telling whether a banana or a fish is healthy

#### Theoretical numerical side

- A test-box for generic optimization programs : NMF is a constrained non-convex (but biconvex) problem
- Robustness analysis of algorithm
- Tensor
- Sparsity

#### Analytical side

Non-negative rank rank<sup>+</sup> := smallest r such that

$$\mathbf{X} = \sum_{i=1}^r \mathbf{X}_i, \quad : \; \mathbf{X}_i \;$$
 rank-1 and non-negative.

How to find / estimate / bound rank^+, e.g.  $\mathsf{rank}_{\mathsf{psd}}(\mathbf{X}) \leq \mathsf{rank}^+(\mathbf{X}).$ 

- Extended formulations and combinatorics
- Log-rank Conjecture of communication system
- 3-SAT, Exponential time hypothesis,  $\mathbf{P} \neq \mathbf{NP}$

#### $15 \, / \, 119$

# Outline

#### Introduction

## 2 Find $(\mathbf{W},\mathbf{H})$ numerically

- Variations on BCD
- A-HALS
- Projected Gradient Update and the Multiplicative update

#### ${f 3}$ Find $({f W},{f H})$ numerically fast : acceleration via extrapolation

- Recall : acceleration in single variable problem
- Accelerating NMF algorithms using extrapolation
- Convergence of the algorithms
  - Application of PALM on NMF
  - Convergence condition of PALM

Problem  $(\mathcal{P})$  : Given  $(\mathbf{X}, r)$ , solve

$$[\mathbf{W}, \ \mathbf{H}] = \underset{\mathbf{W} \ge \mathbf{0}, \mathbf{H} \ge \mathbf{0}}{\operatorname{arg\,min}} \Phi(\mathbf{W}, \mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}.$$

- Equivalent objective function :  $rac{1}{2} \| \mathbf{X} \mathbf{W} \mathbf{H} \|_F^2$ .
- Simplify notation : hide some  $\geq \mathbf{0}, \frac{1}{2}, F$  and just write

$$\min_{\mathbf{W},\mathbf{H}} \Phi(\mathbf{W},\mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|^2.$$

# Standard framework to solve $(\mathcal{P})$

Problem  $(\mathcal{P})$  :  $\min_{\mathbf{W},\mathbf{H}} \Phi(\mathbf{W},\mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|^2$ . Approach : BCD (Block Coordinate Descent)<sup>2</sup>

Algorithm BCD framework for  ${\cal P}$ 

Input:  $X \in \mathbb{R}^{m \times n}_+$ ,  $r \in \mathbb{N}$ , an initialization  $W \in \mathbb{R}^{m \times r}_+$ ,  $H \in \mathbb{R}^{r \times n}_+$ Output: W and H

1: for 
$$k = 1, 2, ...$$
 do

- 2: Update[ $\mathbf{W}$ ] : do something with  $\Phi, \mathbf{X}, \mathbf{W}, \mathbf{H}$ .
- 3: Update[H] : do something with  $\Phi$ , X, W, H.

4: end for

The goal of "do something" is to achieve

$$\Phi(\mathbf{W}^{k+1}, \mathbf{H}^{k+1}) \le \Phi(\mathbf{W}^{k+1}, \mathbf{H}^k) \le \Phi(\mathbf{W}^k, \mathbf{H}^k).$$

 $^2$ Other names : Gauss-Seidel iteration, alternating minimization (for 2 blocks)  $18 \ / \ 119$ 

# An example

#### Algorithm BCD framework for $\mathcal{P}$

**Input:**  $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ ,  $r \in \mathbb{N}$ , an initialization  $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ ,  $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ **Output:**  $\mathbf{W}$  and  $\mathbf{H}$ 

- 1: for k = 1, 2, ... do 2: Update[W] as  $\mathbf{W} \leftarrow \arg\min_{\mathbf{W} \ge 0} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$ . 3: Update[H] as  $\mathbf{H} \leftarrow \arg\min_{\mathbf{H} \ge 0} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$ .
- 4: end for

# An example

#### Algorithm BCD framework for $\mathcal{P}$

Input:  $X \in \mathbb{R}^{m \times n}_+$ ,  $r \in \mathbb{N}$ , an initialization  $W \in \mathbb{R}^{m \times r}_+$ ,  $H \in \mathbb{R}^{r \times n}_+$ Output: W and H

1: for k = 1, 2, ... do 2: Update[W] as  $\mathbf{W} \leftarrow \arg\min_{\mathbf{W} \ge 0} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$ . 3: Update[H] as  $\mathbf{H} \leftarrow \arg\min_{\mathbf{H} \ge 0} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$ .

4: end for

Symmetry :  $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 = \|\mathbf{X}^\top - \mathbf{H}^\top\mathbf{W}^\top\|_F^2$ ,  $\rightarrow$  we can use the same scheme on both variables. We can focus on one variable, says  $\mathbf{H}$  (or  $\mathbf{W}$ ).

If asymmetric regularization exists on  ${f W}$  (or  ${f H}$ ) : we have to handle them separately. 20/119

Block partitions : on how coordinate is being defined<sup>†</sup>.
 This talk : coordinate is H (matrix) or H(i, :) (vector).

- Block partitions : on how coordinate is being defined<sup>†</sup>.
   This talk : coordinate is H (matrix) or H(i, :) (vector).
- Index selection (indexing) : on how coordinate is being selected<sup>#</sup>. This talk : cyclic indexing and A-HALS.

- Block partitions : on how coordinate is being defined<sup>†</sup>.
   This talk : coordinate is H (matrix) or H(i, :) (vector).
- Index selection (indexing) : on how coordinate is being selected<sup>#</sup>. This talk : cyclic indexing and A-HALS.
- Update scheme : on how coordinate is being updated<sup>#</sup>. This talk : "exact" coordinate minimization using 1st order method (e.g. gradient descent). Exact = working on the original objective function, no modification. Inexact = working on modified objective function. e.g. consider relaxation.

- Block partitions : on how coordinate is being defined<sup>†</sup>.
   This talk : coordinate is H (matrix) or H(i, :) (vector).
- Index selection (indexing) : on how coordinate is being selected<sup>#</sup>. This talk : cyclic indexing and A-HALS.
- Update scheme : on how coordinate is being updated<sup>#</sup>. This talk : "exact" coordinate minimization using 1st order method (e.g. gradient descent). Exact = working on the original objective function, no modification. Inexact = working on modified objective function. e.g. consider relaxation.
- Other variants (not in this talk)

 $\dagger$  Kim-He-Park 2014," Algo. for nonnegative matrix and tensor factorizations: a unified view based on block coordinate descent framework" J. Global Optimization.

#Shi-Tu-Xu-Yin 2017," A primer on coordinate descent algorithms." arXiv:1610.00040

# The idea of HALS and A-HALS

Says coordinates are vectors (col. of W and row of H), we have  $\Phi = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr} \langle \mathbf{X}_i, \mathbf{w}_i \mathbf{h}_i \rangle + c.$ 

# The idea of HALS and A-HALS

Says coordinates are vectors (col. of W and row of H), we have  $\Phi = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr} \langle \mathbf{X}_i, \mathbf{w}_i \mathbf{h}_i \rangle + c.$ 

#### Alternating minimization using cyclic indexing Other name : BCD with r = 2 with cyclic component selection Domain name in NMF : HALS (Hierarchical alternating least squares<sup>†</sup>)

Update order :  $\mathbf{w}_1 \rightarrow \mathbf{h}_1 \rightarrow \mathbf{w}_2 \rightarrow \mathbf{h}_2 \rightarrow \mathbf{w}_3 \rightarrow \mathbf{h}_3 \rightarrow ...$ 

# The idea of HALS and A-HALS

Says coordinates are vectors (col. of W and row of H), we have  $\Phi = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr} \langle \mathbf{X}_i, \mathbf{w}_i \mathbf{h}_i \rangle + c.$ 

### Alternating minimization using cyclic indexing Other name : BCD with r = 2 with cyclic component selection Domain name in NMF : HALS (Hierarchical alternating least squares<sup>†</sup>)

Update order :  $\mathbf{w}_1 \rightarrow \mathbf{h}_1 \rightarrow \mathbf{w}_2 \rightarrow \mathbf{h}_2 \rightarrow \mathbf{w}_3 \rightarrow \mathbf{h}_3 \rightarrow ...$ 

**A-HALS**<sup>#</sup> (Accelerated-HALS) A special kinds of cyclic coordinate selection

$$\begin{array}{c} \mathsf{U}\mathsf{p}\mathsf{d}\mathsf{a}\mathsf{t}\mathsf{e} \ \mathsf{order}: \ \underbrace{\mathbf{w}_1 \to \mathbf{w}_2 \to \cdots \to \mathbf{w}_r}_{\mathsf{several times}!!} \to \underbrace{\mathbf{h}_1 \to \mathbf{h}_2 \to \cdots \to \mathbf{h}_r}_{\mathsf{several times}!!} \to \dots \end{array}$$

† Cichocki-Zdunke-Amari 2007, "Hierarchical ALS Algorithms for Nonnegative Matrix and 3D Tensor Factorization", International Conf. on ICA.

# Gillis-Glineur 2012, "Accelerated Multiplicative Updates and Hierarchical ALS Algo. for NMF", Neural Computation.  $$27\,/\,119$$ 

## A-HALS = avoids repeated computations + re-uses

Projected<sup>†</sup> gradient descent with step size  $t \ge 0$ 

$$\mathbf{w}_{i} = \mathbf{w}_{i} - t \underbrace{(\|\mathbf{h}_{i}\|_{2}^{2}\mathbf{w}_{i} - \mathbf{X}_{i}\mathbf{h}_{i}^{\top})}_{\nabla_{\mathbf{w}_{i}}\Phi}, \quad \mathbf{h}_{i} = \mathbf{h}_{i} - t \underbrace{(\|\mathbf{w}_{i}\|_{2}^{2}\mathbf{h}_{i} - \mathbf{w}_{i}^{\top}\mathbf{X})}_{\nabla_{\mathbf{h}_{i}}\Phi}.$$

## A-HALS = avoids repeated computations + re-uses

Projected<sup>†</sup> gradient descent with step size  $t \ge 0$  $\mathbf{w}_i = \mathbf{w}_i - t(\|\mathbf{h}_i\|_2^2 \mathbf{w}_i - \mathbf{X}_i \mathbf{h}_i^{\top}), \quad \mathbf{h}_i = \mathbf{h}_i - t(\|\mathbf{w}_i\|_2^2 \mathbf{h}_i - \mathbf{w}_i^{\top} \mathbf{X}).$  $\nabla_{\mathbf{w}_i} \Phi$  $\nabla_{\mathbf{h}_i} \Phi$ Algorithm A-HALS **Algorithm** HALS 1:  $\mathbf{w}_1 = \mathbf{w}_1 - t(\|\mathbf{h}_1\|_2^2 \mathbf{w}_1 - \mathbf{X}_1 \mathbf{h}_1^\top)$ 1: Compute  $\mathbf{A} = \mathbf{H}\mathbf{H}^{\top}$ ,  $\mathbf{B} = \mathbf{X}\mathbf{H}^{\top}$ 2:  $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ 2:  $\mathbf{w}_1 = \mathbf{w}_1 - t(\|\mathbf{h}_1\|_2^2 \mathbf{w}_1 - \mathbf{X}_1 \mathbf{h}_1^{\top})$ **3**:  $\mathbf{w}_2 = \mathbf{w}_2 - t(\|\mathbf{h}_2\|_2^2 \mathbf{w}_2 - \mathbf{X}_2 \mathbf{h}_2^\top)$ **3**:  $\mathbf{w}_2 = \mathbf{w}_2 - t(\|\mathbf{h}_2\|_2^2 \mathbf{w}_2 - \mathbf{X}_2 \mathbf{h}_2^\top)$ **4**:  $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ 4:  $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_3^{\top})$ 5: Compute  $\mathbf{C} = \mathbf{W}^\top \mathbf{W}, \mathbf{D} = \mathbf{W}^\top \mathbf{X}$ **5**:  $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_2^{\top})$ **6**:  $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ **6**:  $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_2^\top \mathbf{X}_3)$ 7:  $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ 7: ... 8:  $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_3^\top \mathbf{X}_3)$ 9:

A-HALS : Line 2-4, 6-8 repeated a few times.

# A-HALS = avoids repeated computations + re-uses

Projected<sup>†</sup> gradient descent with step size  $t \ge 0$  $\mathbf{w}_i = \mathbf{w}_i - t \left( \|\mathbf{h}_i\|_2^2 \mathbf{w}_i - \mathbf{X}_i \mathbf{h}_i^\top \right), \quad \mathbf{h}_i = \mathbf{h}_i - t \left( \|\mathbf{w}_i\|_2^2 \mathbf{h}_i - \mathbf{w}_i^\top \mathbf{X} \right).$  $\nabla_{\mathbf{w}} \Phi$  $\nabla_{\mathbf{h}_i} \Phi$ Algorithm HALS Algorithm A-HALS 1: Compute  $\mathbf{A} = \mathbf{H}\mathbf{H}^{\top}$ ,  $\mathbf{B} = \mathbf{X}\mathbf{H}^{\top}$ 1:  $\mathbf{w}_1 = \mathbf{w}_1 - t(\|\mathbf{h}_1\|_2^2 \mathbf{w}_1 - \mathbf{X}_1 \mathbf{h}_1^\top)$ 2:  $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ 2:  $\mathbf{w}_1 = \mathbf{w}_1 - t(\|\mathbf{h}_1\|_2^2 \mathbf{w}_1 - \mathbf{X}_1 \mathbf{h}_1^{\top})$ **3**:  $\mathbf{w}_2 = \mathbf{w}_2 - t(\|\mathbf{h}_2\|_2^2 \mathbf{w}_2 - \mathbf{X}_2 \mathbf{h}_2^\top)$ **3**:  $\mathbf{w}_2 = \mathbf{w}_2 - t(\|\mathbf{h}_2\|_2^2 \mathbf{w}_2 - \mathbf{X}_2 \mathbf{h}_2^\top)$ 4:  $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_3^{\top})$ **4**:  $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ 5: Compute  $\mathbf{C} = \mathbf{W}^\top \mathbf{W}, \mathbf{D} = \mathbf{W}^\top \mathbf{X}$ **5**:  $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_2^{\top})$ **6**:  $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ **6**:  $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_2^\top \mathbf{X}_3)$ 7:  $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ 7: ... 8:  $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_3^\top \mathbf{X}_3)$ 9: ...

A-HALS : Line 2-4, 6-8 repeated a few times. A-HALS avoids repeated computations of *constant terms* :

$$\mathbf{H}\mathbf{H}_{(2n-1)m^2}^{\top}, \ \mathbf{X}\mathbf{H}_{(2n-1)mr}^{\top}, \ \mathbf{W}^{\top}\mathbf{W}_{(2r-1)m^2}, \ \mathbf{W}^{\top}\mathbf{X}_{(2m-1)rn},$$

30 / 119

pre-computing and re-use of these terms gain extra efficiency, improvement is significant for big data<sup>#"</sup> — always A-HALS!

# The projected gradient descent update

The Projected Gradient Descent update of  ${\bf W}$  :

$$\mathbf{W}^{k+1} = \operatorname{Proj}_{\mathbb{R}_+} \Big( \mathbf{W}^k - t \nabla \Phi(\mathbf{W}^k, \mathbf{H}) \Big).$$

How to pick the step-size ?

# The projected gradient descent update

The Projected Gradient Descent update of W :

$$\mathbf{W}^{k+1} = \operatorname{Proj}_{\mathbb{R}_+} \Big( \mathbf{W}^k - t \nabla \Phi(\mathbf{W}^k, \mathbf{H}) \Big).$$

How to pick the step-size ? A simple scheme  $t = \frac{1}{L_{\Phi_{\mathbf{W}}}}$ .

In words : pick step-size as  $L_{\Phi_{\mathbf{W}}}^{-1}$ , where  $L_{\Phi_{\mathbf{W}}}$  = the Lipschitz constant of  $\nabla_{\mathbf{W}} \Phi$  (smoothness constant).

# The projected gradient descent update

The Projected Gradient Descent update of W :

$$\mathbf{W}^{k+1} = \operatorname{Proj}_{\mathbb{R}_+} \Big( \mathbf{W}^k - t \nabla \Phi(\mathbf{W}^k, \mathbf{H}) \Big).$$

How to pick the step-size ? A simple scheme  $t = \frac{1}{L_{\Phi_{\mathbf{W}}}}$ .

In words : pick step-size as  $L_{\Phi_{\mathbf{W}}}^{-1}$ , where  $L_{\Phi_{\mathbf{W}}}$  = the Lipschitz constant of  $\nabla_{\mathbf{W}} \Phi$  (smoothness constant).

$$L_{\Phi_{\mathbf{W}}} = \text{largest singular value of } \mathbf{H}\mathbf{H}^{\top}$$
  
 $\operatorname{Proj}_{\mathbb{R}_{+}}$  is basically  $[ \cdot ]_{+} = \max\{\cdot, 0\}$ .  
Hence in close form :

$$\mathbf{W}^{k+1} = \left[\mathbf{W}^k - \frac{1}{\sigma_{\max}(\mathbf{H}\mathbf{H}^{\top})} \nabla \Phi(\mathbf{W}^k, \mathbf{H})\right]_+$$

PGD update is much faster than the *Multiplicative Update* (MU).

MU :

• It takes a small step size t such that  $\mathbf{W}^{k+1}$  stays within  $\mathrm{I\!R}_+\text{, no projection.}$ 

$$\mathbf{W}^{k+1} = \mathbf{W}. * rac{\mathbf{X}\mathbf{H}^{ op}}{\mathbf{W}^k\mathbf{H}\mathbf{H}^{ op}},$$

where \* is Hadamard product and the division is Hadamard quotient.

• It converges **very slowly**. In general, don't use MU. Pourquoi/why: to make sure  $\mathbf{W}$  stays within  $\mathrm{IR}_+$ , MU take small step  $\implies$  slow !

PGD :

- $\bullet$  It takes reasonably large step size, and IF moved outside  $\rm I\!R_+$  THEN project back.
- $\operatorname{Proj}_{\mathbb{R}_+}$  practically costs nothing unless the data size is  $10^{86}$ .



 $\begin{array}{l} \mathsf{MU}=\mathsf{timid}, \,\mathsf{shy} \,\, \mathsf{person} \,\, \mathsf{that} \,\, \mathsf{is} \,\, \mathsf{too} \,\, \mathsf{cautious} \,\, \mathsf{on} \,\, \mathsf{making} \,\, \mathsf{mistake}.\\ \mathsf{PGD}=\mathsf{brave} \,\, \mathsf{person} \,\, \mathsf{that} \,\, \mathsf{is} \,\, \mathsf{fine} \,\, \mathsf{of} \,\, \mathsf{making} \,\, \mathsf{mistake} \,\, \mathsf{by} \,\, \mathsf{doing} \,\, \mathsf{correction}.\\ \mathsf{Here} \,\, "\,\mathsf{mistake}" \,\, = \, "\,\mathsf{outside} \,\, \mathrm{IR}_+" \,, \, "\,\mathsf{correction}" \,\, = \, "\,\mathrm{Proj}_{\mathrm{IR}_+} \,\, \overset{"}{_{35}} \,/ \, 119 \end{array}$ 

# Outline



## Find $(\mathbf{W}, \mathbf{H})$ numerical

- Variations on BCD
- A-HALS
- Projected Gradient Update and the Multiplicative update

## $\mathbf{3}$ Find $(\mathbf{W},\mathbf{H})$ numerically fast : acceleration via extrapolation

- Recall : acceleration in single variable problem
- Accelerating NMF algorithms using extrapolation
- Convergence of the algorithms
  - Application of PALM on NMF
  - Convergence condition of PALM
#### Let's accelerate !



What's the acceleration for : obtain a local solution  $faster_{37/119}$ 

### Recall : acceleration in single variable problem

 $\text{Problem } \min_{x \in \mathcal{C}} f(x) \text{, } \mathcal{C} \text{ convex set.}$ 

#### Recall : acceleration in single variable problem

```
Problem \min_{x \in C} f(x), C convex set.
At step k:
No acceleration : x_{k+1} = \text{Update}[x_k].
With acceleration : x_{k+1} = \text{Update}[y_k], y_{k+1} = \text{Extrapolate}[x_{k+1}, x_k].
```

Problem  $\min_{x \in \mathcal{C}} f(x)$ ,  $\mathcal{C}$  convex set. At step k: No acceleration :  $x_{k+1} = Update[x_k]$ . With acceleration :  $x_{k+1} = \text{Update}[y_k], y_{k+1} = \text{Extrapolate}[x_{k+1}, x_k].$ To be specific : PGD Update  $x_{k+1} = \operatorname{Proj}_{\mathcal{C}}(x_k - t_k \nabla f(x_k)).$ Linear extrapolation  $x_{k+1} = \operatorname{Proj}_{\mathcal{C}}(y_k - t_k \nabla f(y_k)).$  $y_{k+1} = x_{k+1} + \beta_k (x_{k+1} - x_k).$ 

i.e. Extrapolate $[x_{k+1}, x_k]$  is modeled by  $\beta_k$ : a single extrapolation parameter.

#### Why extrapolation : gradient descent zig-zags on ellipse

Facts : consecutive update directions of GD are orthogonal ( $\perp$ ). If the landscape is not "spherical", GD zig-zags  $\rightarrow$  slow.

e.g. : moving along a long narrow valley.



Picture modified from http://www.nbertagnolli.com/jekyll/update/2015/10/28/Descent-Methods.html

### An slide from my other slides



Picture from https://angms.science/doc/teaching/GDLS.pdf

### What machine learning people do to counter zig-zag?

**Do tricks on step size** : don't move with step size t but  $\frac{\iota}{\text{damping factor}}$ 



The idea behind **AdaGrad** and **AdaDelta** : shrink the step size when you see zig-zag (trace of the objective function appears to plateau).

Do tricks on direction : by extrapolation with momentum.



 $\label{eq:ldea:apply} \begin{array}{ll} \mbox{Idea}: \mbox{ apply extrapolation}.\\ \mbox{Extrapolate} = \mbox{add gradient history}. \end{array}$ 

(1) if gradients in consecutive steps have consistent direction

 $\implies$  extrapolate = accelerate.

(2) if gradients in consecutive steps oscillates (continuously changing direction)

 $\implies$  extrapolate = damp oscillation = acceleration.

Figure shows the trace of points decomposed into x- and y-component. The x-components have consistent direction while y-components are not.

 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$ 



 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k (x_{k+1} - x_k).$ 



 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k (x_{k+1} - x_k).$ 



 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$ 



 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$ 



 $x_{k+1} = \mathsf{Update}[y_k], \ y_{k+1} = x_{k+1} + \beta_k (x_{k+1} - x_k).$ 



We always have

 $\angle (x_{k+1} - y_k) \ge \angle (x_{k+2} - x_{k+1}) \ge \angle (x_{k+2} - y_{k+1})$ 

i.e. the direction of the last step is **in between** the directions of previous two gradient steps : zig-zag effect is reduced !



For convex function,

$$\beta_k = \frac{1 - \alpha_k}{\alpha_{k+1}}, \ \alpha_{k+1} = \frac{1 + \sqrt{1 + 4\alpha_k^2}}{2}, \alpha_1 \in (0, 1)$$

**2** For **smooth strongly convex** function with *conditional number* Q,

$$\beta_k = \frac{1-\sqrt{Q}}{1+\sqrt{Q}}, \text{ where } Q = \frac{L}{\mu} = \frac{\text{Smoothness parameter}}{\text{Strong convexity parameter}}$$

With convergence improvement : from  $\mathcal{O}(Q \log \frac{1}{\epsilon})$  to  $\mathcal{O}(\sqrt{Q} \log \frac{1}{\epsilon})$ 

Key : Nesterov's acceleration has a close-form formula for  $\beta_k$ 

#### A slide from my other slides

#### Other $\beta_k$ schemes

Nesterov's parameters looks so complicated

$$\alpha_{k+1} = \frac{\sqrt{\alpha_k^4 + 4\alpha_k^2 - \alpha_k^2}}{2}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

Another Nesterov's parameters

$$\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + \kappa^{-1}\alpha_{k+1}, \ \beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}$$

Yet another Nesterov's parameters

$$\alpha_{k+1} = \frac{1 + \sqrt{1 + 4\alpha_k^2}}{2}, \ \beta_k = \frac{1 - \alpha_k}{\alpha_{k+1}}.$$

Paul Tseng parameter

$$\beta_k = \frac{k-1}{k+2}.$$

Using conditional number

$$\beta_{k} = \beta = \frac{1 - \sqrt{\kappa'}}{1 + \sqrt{\kappa'}}, \quad \kappa' = \frac{1}{\kappa}, \quad \kappa = \frac{\sigma_{\max}(\mathbf{Q})}{\sigma_{\min}(\mathbf{Q})} = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$$

$$70 / 14$$

Picture from https://angms.science/doc/teaching/GDLS.pdf

Key : Nesterov's acceleration has a close-form formula for  $\beta_k$ 

#### Extrapolation is not monotone, nor descent, nor greedy

GD is locally optimal/greedy ⇒ extrapolation may ↑objective value • Extrapolation = a risky move



Picture from Donoghue-Candés 2015, "Adaptive Restart for Accelerated Gradient Schemes" Acceleration comes from doing the risky move :

"sacrifice the decreases of objective value now for the better future"

Actually also sacrifice robustness : accelerated gradient is not stable to noise (Devolder-Glineur-Nesterov 2014) 54/119

### A slide from my other slides



Picture from https://angms.science/doc/teaching/GDLS.pdf

#### Our case : NMF is not cvx

Problem ( $\mathcal{P}$ ) : {Given ( $\mathbf{X}, r$ ), solve  $\min_{\mathbf{W}, \mathbf{H}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|^2, \mathbf{W}, \mathbf{H} \in \mathbb{R}_+$ } is non-cvx but bi-cvx.

 $\implies$  no strong cvx parameter  $\mu$ . Cannot use expression likes  $\beta_k = \frac{1 - \sqrt{Q}}{1 + \sqrt{Q}}$ .

#### Our case : NMF is not cvx

Problem  $(\mathcal{P})$ : {Given  $(\mathbf{X}, r)$ , solve  $\min_{\mathbf{W}, \mathbf{H}} ||\mathbf{X} - \mathbf{W}\mathbf{H}||^2, \mathbf{W}, \mathbf{H} \in \mathbb{R}_+$ } is non-cvx but bi-cvx.

 $\implies$  no strong cvx parameter  $\mu$ . Cannot use expression likes  $\beta_k = \frac{1 - \sqrt{Q}}{1 + \sqrt{Q}}$ . For the acceleration scheme of the two variables

 $\left\{ \begin{array}{l} {\rm On}\; \mathbf{W} & \left\{ \begin{array}{l} {\rm Update} \;\; \mathbf{W}_{\sf new} = {\rm Update}[\mathbf{Y}_{\sf old}, \mathbf{H}_{\sf old}] \\ \\ {\rm Extrapolate} \;\; \mathbf{Y}_{\sf new} = \mathbf{W}_{\sf new} + \beta_k^{\mathbf{W}}(\mathbf{W}_{\sf new} - \mathbf{W}_{\sf old}) \\ \\ {\rm On}\; \mathbf{H} \;\; \left\{ \begin{array}{l} {\rm Update} \;\; \mathbf{H}_{\sf new} = {\rm Update}[\mathbf{W}_{\sf new}, \mathbf{G}_{\sf old}] \\ \\ {\rm Extrapolate} \;\; \mathbf{G}_{\sf new} = \mathbf{H}_{\sf new} + \beta_k^{\mathbf{H}}(\mathbf{H}_{\sf new} - \mathbf{H}_{\sf old}) \end{array} \right. \end{array} \right.$ 

Need a way (close-/no close-form) to find  $\beta_k$  !

#### Our case : NMF is not cvx

Problem  $(\mathcal{P})$ : {Given  $(\mathbf{X}, r)$ , solve  $\min_{\mathbf{W}, \mathbf{H}} ||\mathbf{X} - \mathbf{W}\mathbf{H}||^2, \mathbf{W}, \mathbf{H} \in \mathbb{R}_+$ } is non-cvx but bi-cvx.

 $\implies$  no strong cvx parameter  $\mu$ . Cannot use expression likes  $\beta_k = \frac{1 - \sqrt{Q}}{1 + \sqrt{Q}}$ . For the acceleration scheme of the two variables

 $\left\{ \begin{array}{l} {\rm On}\; \mathbf{W} \; \left\{ \begin{array}{l} {\rm Update} \;\; \mathbf{W}_{{\rm new}} = {\rm Update}[\mathbf{Y}_{{\rm old}}, \mathbf{H}_{{\rm old}}] \\ \\ {\rm Extrapolate} \;\; \mathbf{Y}_{{\rm new}} = \mathbf{W}_{{\rm new}} + \beta_k^{\mathbf{W}}(\mathbf{W}_{{\rm new}} - \mathbf{W}_{{\rm old}}) \\ \\ {\rm On}\; \mathbf{H} \; \left\{ \begin{array}{l} {\rm Update} \;\; \mathbf{H}_{{\rm new}} = {\rm Update}[\mathbf{W}_{{\rm new}}, \mathbf{G}_{{\rm old}}] \\ \\ {\rm Extrapolate} \;\; \mathbf{G}_{{\rm new}} = \mathbf{H}_{{\rm new}} + \beta_k^{\mathbf{H}}(\mathbf{H}_{{\rm new}} - \mathbf{H}_{{\rm old}}) \end{array} \right. \end{array} \right.$ 

Need a way (close-/no close-form) to find  $\beta_k$  !

Approach : an ad hoc heurisitic in the "line search" style.  $\frac{58}{119}$ 

Why ad hoc heuristics ?

- (1) The ncvx problem is hard.
- (2) No better idea.
- No convergence theorem now yet (because of (1)).

What's so good ?

- Just a parameter tuning problem.
- Easy to implement.
- Easy to extend to other models.
- Faster than state-of-the-art methods with theoretical convergence proof !

† Xu-Yin 2013 "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion". SIAM J. Img Sci.

- The key  $\beta_k$ 
  - $\beta$  has to be smaller than 1 (same as the convex case)

60 /

• If  $\beta \in (0,1)$  : extrapolation, doing risky step

The key  $\beta_k$ 

- $\beta$  has to be smaller than 1 (same as the convex case)
- If  $\beta \in (0,1)$  : extrapolation, doing risky step
- If  $\beta = \{1,0\}$  : doing {very risky, no} extrapolation

The key  $\beta_k$ 

- $\beta$  has to be smaller than 1 (same as the convex case)
- $\bullet$  If  $\beta \in (0,1)$  : extrapolation, doing risky step
- If  $\beta = \{1,0\}$  : doing {very risky, no} extrapolation
- Can't use line search<sup>†</sup> to find  $\beta$  : experimentally found  $\beta$  close to 0

- minor extrapolation, effectively doing nothing

The key  $\beta_k$ 

- $\beta$  has to be smaller than 1 (same as the convex case)
- $\bullet$  If  $\beta \in (0,1)$  : extrapolation, doing risky step
- If  $\beta = \{1,0\}$  : doing {very risky, no} extrapolation
- Can't use line search<sup>†</sup> to find  $\beta$  : experimentally found  $\beta$  close to 0

- minor extrapolation, effectively doing nothing

## In the "walking person metaphor" :

- MU shy person walking in caution with small step size
- PGD brave person walking in reasonably step size
- E-PGD ambious person walking in big step size

# Details : Update[ $\beta_k$ ]

Landscape of variable at each iteration is different  $\implies$  dynamical update

Algorithm A dynamic line search style<sup>†</sup> ad hoc heuristics

**Input:** Parameters  $1 < \bar{\gamma} < \gamma < \eta$ , an initialization  $\beta_1 \in (0, 1)$ **Output:**  $\beta_k$ : the extrapolation parameter

- 1: Set  $\bar{\beta} = 1$  (dynamic "upper bound" of  $\beta$ )
- 2: if error  $\downarrow$  at iteration k then
- 3: Increase  $\beta_{k+1} : \beta_{k+1} = \min\{\bar{\beta}, \gamma\beta_k\}$
- 4: (Increase  $\bar{\beta}$  if  $\bar{\beta} < 1$  :  $\bar{\beta} = \min\{1, \bar{\gamma}\bar{\beta}\}$ )
- 5: **else**
- 6: Decrease  $\beta_{k+1}$  :  $\beta_{k+1} = \beta_k/\eta$
- 7: Set  $\bar{\beta} = \beta_k$
- 8: end if

 $\gamma\text{, }\bar{\gamma}\text{, }\eta$  : growth and decay parameters

 $\dagger$ Line search after updates of W and H – performed after the update!

### Meaning

- Go further/"speed up" when suitable (error  $\downarrow$ ) : more ambitious, make  $\beta \uparrow$ , take more risk

- Go back/"slow down" when not suitable (error $\uparrow$ ) : less ambitious, make  $\beta \downarrow$ , take less risk



#### The full algo of Accelerated NMF using extrapolation

**Input:** X, initialization W, H, parameters  $hp \in \{1, 2, 3\}$  (extrapolation/projection of H). Output: W.H. 1:  $\mathbf{W}_{u} = \mathbf{W}; \mathbf{H}_{u} = \mathbf{H}; e(0) = ||\mathbf{X} - \mathbf{W}\mathbf{H}||_{F}$ . 2: for  $k = 1, 2, \ldots$  do 3: Compute  $\mathbf{H}_n$  by  $\min_{\mathbf{H}_n \geq 0} ||\mathbf{X} - \mathbf{W}_y \mathbf{H}_n||_F^2$  using  $\mathbf{H}_y$  as initial iterate. 4: 5: 6: 7: 8: 9: 10: if hp > 2 then Extrapolate:  $\mathbf{H}_{u} = \mathbf{H}_{n} + \beta_{k}(\mathbf{H}_{n} - \mathbf{H}).$ end if if hp = 3 then Project:  $\mathbf{H}_{y} = \max(0, \mathbf{H}_{y}).$ end if Compute  $\mathbf{W}_n$  by  $\min_{\mathbf{W}_n > 0} ||\mathbf{X} - \mathbf{W}_n \mathbf{H}_y||_F^2$  using  $\mathbf{W}_y$  as initial iterate. 11: Extrapolate:  $\mathbf{W}_{u} = \mathbf{W}_{n} + \beta_{k}(\mathbf{W}_{n} - \mathbf{W}).$ 12: if hp = 1 then 13: Extrapolate:  $\mathbf{H}_{u} = \mathbf{H}_{n} + \beta_{k}(\mathbf{H}_{n} - \mathbf{H}).$ 14: end if 15: Compute error:  $e(k) = ||\mathbf{X} - \mathbf{W}_n \mathbf{H}_u||_F$ . 16: if e(k) > e(k-1) then 17: Restart:  $\mathbf{H}_{u} = \mathbf{H}_{n}$ ;  $\mathbf{W}_{u} = \mathbf{W}_{n}$ . 18: 19: 20: else  $\mathbf{H} = \mathbf{H}_n$ :  $\mathbf{W} = \mathbf{W}_n$ . end if 21: end for

Notation :  $\mathbf{W}_n$  normal variable,  $\mathbf{W}_y$  extrpolate variable,  $\mathbf{W}$  previous  $\mathbf{W}_n$  ... too hard to read !!

# Algorithm (hp = 1), simplified

Input:  $\mathbf{X}$ , initialization  $\mathbf{W}$ ,  $\mathbf{H}$ Output:  $\mathbf{W}$ ,  $\mathbf{H}$ 

1: 
$$\mathbf{W}_y = \mathbf{W}$$
;  $\mathbf{H}_y = \mathbf{H}$ ;  $e(0) = ||\mathbf{X} - \mathbf{W}\mathbf{H}||_F$ .

2: for 
$$k = 1, 2, ...$$
 do

- 3: **Up**date[ $\mathbf{H}_n$ ] w.r.t.  $\mathbf{H}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$  using  $\mathbf{H}_y$  as initial iterate.
- 4: **Up**date[ $\mathbf{W}_n$ ] wr.t.  $\mathbf{W}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$  using  $\mathbf{W}_y$  as initial iterate.
- 5: **Extrapolate**[ $\mathbf{W}_{y}$ ] :  $\mathbf{W}_{y} = \mathbf{W}_{n} + \beta_{k}(\mathbf{W}_{n} \mathbf{W})$ .
- 6: **Ex**trapolate[ $\mathbf{H}_y$ ] :  $\mathbf{H}_y = \mathbf{H}_n + \beta_k(\mathbf{H}_n \mathbf{H})$ .
- 7: Compute error:  $e(k) = ||\mathbf{X} \mathbf{W}_n \mathbf{H}_y||_F$ .
- 8: if e(k) > e(k-1) then

9: Restart: 
$$\mathbf{H}_y = \mathbf{H}_n$$
;  $\mathbf{W}_y = \mathbf{W}_n$ 

10: else

11: 
$$\mathbf{H} = \mathbf{H}_n; \ \mathbf{W} = \mathbf{W}_n.$$

- 12: end if
- 13: end for

# "Up, Up, Ex, Ex"

# Algorithm (hp = 2), simplified

Input:  $\mathbf{X}$ , initialization  $\mathbf{W}$ ,  $\mathbf{H}$ Output:  $\mathbf{W}$ ,  $\mathbf{H}$ 

- 1:  $\mathbf{W}_y = \mathbf{W}$ ;  $\mathbf{H}_y = \mathbf{H}$ ;  $e(0) = ||\mathbf{X} \mathbf{W}\mathbf{H}||_F$ .
- 2: for k = 1, 2, ... do
- 3: **Up**date[ $\mathbf{H}_n$ ] w.r.t.  $\mathbf{H}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$  using  $\mathbf{H}_y$  as initial iterate.
- 4: **Extrapolate** $[\mathbf{H}_y]$ :  $\mathbf{H}_y = \mathbf{H}_n + \beta_k (\mathbf{H}_n \mathbf{H}).$
- 5: **Up**date[ $\mathbf{W}_n$ ] wr.t.  $\mathbf{W}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$  using  $\mathbf{W}_y$  as initial iterate.
- 6: **Extrapolate** $[\mathbf{W}_y]$ :  $\mathbf{W}_y = \mathbf{W}_n + \beta_k (\mathbf{W}_n \mathbf{W}).$
- 7: Compute error:  $e(k) = ||\mathbf{X} \mathbf{W}_n \mathbf{H}_y||_F$ .
- 8: if e(k) > e(k-1) then

9: Restart: 
$$\mathbf{H}_y = \mathbf{H}_n$$
;  $\mathbf{W}_y = \mathbf{W}_n$ 

10: else

11: 
$$\mathbf{H} = \mathbf{H}_n; \ \mathbf{W} = \mathbf{W}_n.$$

- 12: end if
- 13: end for

# "Up, Ex, Up, Ex"

# Algorithm (hp = 3), simplified

Input:  $\mathbf{X}$ , initialization  $\mathbf{W}, \mathbf{H}$ Output:  $\mathbf{W}, \mathbf{H}$ 

- 1:  $\mathbf{W}_y = \mathbf{W}$ ;  $\mathbf{H}_y = \mathbf{H}$ ;  $e(0) = ||\mathbf{X} \mathbf{W}\mathbf{H}||_F$ .
- 2: for k = 1, 2, ... do
- 3: **Up**date[ $\mathbf{H}_n$ ] w.r.t.  $\mathbf{H}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$  using  $\mathbf{H}_y$  as initial iterate.
- 4: **Extrapolate** $[\mathbf{H}_y]$ :  $\mathbf{H}_y = \mathbf{H}_n + \beta_k (\mathbf{H}_n \mathbf{H}).$
- 5: **Project**:  $\mathbf{H}_y = \max(0, \mathbf{H}_y)$ .
- 6: **Up**date[ $\mathbf{W}_n$ ] wr.t.  $\mathbf{W}_n \ge 0$  with  $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$  using  $\mathbf{W}_y$  as initial iterate.
- 7: **Extrapolate** $[\mathbf{W}_y]$ :  $\mathbf{W}_y = \mathbf{W}_n + \beta_k (\mathbf{W}_n \mathbf{W}).$
- 8: Compute the error:  $e(k) = ||\mathbf{X} \mathbf{W}_n \mathbf{H}_y||_F$ .
- 9: **if** e(k) > e(k-1) **then**
- 10: Restart:  $\mathbf{H}_y = \mathbf{H}_n$ ;  $\mathbf{W}_y = \mathbf{W}_n$ .
- 11: else
- 12:  $\mathbf{H} = \mathbf{H}; \ \mathbf{W} = \mathbf{W}_n.$
- 13: end if
- 14: end for

"Up, Ex, Pro, Up, Ex"

Extrapolation may break NN (  $\geq 0)$  constraint :

hp = 1		hp = 2		hp = 3	
(Up-Up-Ex-Ex)		(Up-Ex-Up-Ex)		(Up-Ex-Pro-Up-Ex)	
Step	NN?	Step	NN?	Step	NN?
$Update[\mathbf{H}_n]$	Y	$Update[\mathbf{H}_n]$	Y	$Update[\mathbf{H}_n]$	Y
$Update[\mathbf{W}_n]$	Y	$Extrap[\mathbf{H}_{y}]$	Ν	$Extrap[\mathbf{H}_{y}]$	Ν
				$Project[\mathbf{H}_y]$	Y
$Extrap[\mathbf{H}_y]$	N	$Update[\mathbf{W}_n]$	Y	$Update[\mathbf{W}_n]$	Y
$Extrap[\mathbf{W}_y]$	Ν	$Extrap[\mathbf{W}_{y}]$	Ν	$Extrap[\mathbf{W}_{y}]$	Ν

Update using matrix with negative values : Update[ $\mathbf{H}_n$ ] w.r.t.  $\mathbf{H}_n \ge 0$  with  $(\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n)$ , using  $\mathbf{H}_y$  as initial iterate Update[ $\mathbf{W}_n$ ] wr.t.  $\mathbf{W}_n \ge 0$  with  $(\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y)$ , using  $\mathbf{W}_y$  as initial iterate

# Summary and notes (3/3)

Restart using e(k) as  $\|\mathbf{X} - \mathbf{W}_n \mathbf{H}_{y}\|_F$  not  $\|\mathbf{X} - \mathbf{W}_n \mathbf{H}_{n}\|_F$ 

Why :

(i)  $\mathbf{W}_n$  was updated according to  $\mathbf{H}_y$  (see point 2)

(ii) it gives the algorithm some degrees of freedom to possibly increase the objective function  $% \left( {{{\left[ {{{\rm{T}}} \right]}_{{\rm{T}}}}_{{\rm{T}}}} \right)$ 

(iii) computationally cheaper, as compute  $\|\mathbf{X} - \mathbf{W}_n \mathbf{H}_n\|_F$  need O(mnr) operations instead of  $O(mr^2)$  by re-using previous computed terms :

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} = \|\mathbf{X}\|_{F}^{2} - 2\left\langle \mathbf{W}, \mathbf{X}\mathbf{H}^{\top} \right\rangle + \left\langle \mathbf{W}^{\top}\mathbf{W}, \mathbf{H}\mathbf{H}^{\top} \right\rangle$$

Note : if the variables converges, using  $\mathbf{W}_n$ ,  $\mathbf{W}_y$  is effectively the same as in  $\mathbf{W}_n^{\infty} = \mathbf{W}_y^{\infty}$  (after projection)
### Experiments

#### Notations

- A-HALS : vector-wise update, compute approximate solution
- ANLS : subproblem solved exactly using active-set methods
- E : extrapolation

Set up

- Average error over 10 trials
- $\mathbf{W}, \mathbf{H}, \mathbf{X}$  randomly generated  $\sim \mathcal{U}[0, 1]$ , m = n = 200, r = 20
- ullet Real  ${f X}$  from real data is also used.
- Error comparisons : using lowest relative error  $e_{\min}$  across all algorithms, at step k,

$$E(k) = \frac{\|\mathbf{X} - \mathbf{W}^k \mathbf{H}^k\|_F}{\|\mathbf{X}\|_F} - e_{\min}$$

- It is possible  $e_{\min} = 0$  and not shown
- Extrapolation parmater  $\beta_0 = [0.25, 0.5, 0.75]$
- $\eta_0 = [1.5, 2, 3]$
- $\gamma, \bar{\gamma} = [1.01, 1.005], [1.05, 1.01], [1.1, 1.05]$
- For display : only best and worst to illustrate sensitivity (for  $h_{12} \neq p_{2}$ )



# Compare with other method on speed (time)



Average err. of ANLS, A-HALS and extrapolated variants, on low-rank (left) and full-rank (right) synthetic data.

APG-MF<sup> $\dagger$ </sup> = an extrapolated proximal type algorithm, with convergence proof.

#### Fast conclusion : E wins and beats $APG-MF^{\dagger}$ .

 $\dagger$  Xu-Yin 2013 "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion". SIAM J. Img Sci.  $75\ /\ 119$ 

### Overall results : E wins!

Method	Data	Ex wins?
	Low/full rank synthetic data	YES
A-HALS	Dense Image data <sup>†</sup>	YES
	Sparse text data $^{\#}$	YES
	Low/full rank synthetic data	YES
ANLS	Dense Image data <sup>†</sup>	YES
	Sparse text data $^{\#}$	YES

† ORL, Umist, CBCL, Frey.

 $^{\#}$  Zhong-Ghosh 2005. Generative model-based document clustering: a comparative study

#### Conclusions

- No matter what method XXX, E-XXX > XXX.
- E-XXX > APG-MF (an extrapolated proximal-type method).
- Between E-ANLS vs E-A-HALS : no clear winner
  - ▶ Low rank synthetic data : E-ANLS ≫ everything
  - Dense data : E-A-HALS  $\approx$  E-ANLS, although A-HALS > ANLS
  - Sparse data : E-A-HALS  $\gg$  everything
- Between different hp
  - Up-Ex-Up-Ex (hp = 2) seems worst
  - Up-Up-Ex-Ex (hp = 1) or Up-Ex-Pro-Up-Ex (hp = 3) are better

Don't trust me ? Go https://arxiv.org/abs/1805.06604, try the code!

#### A quick-and-dirty test on tensor



# Outline



#### mtroduction

#### 2) Find $(\mathbf{W}, \mathbf{H})$ numerically

- Variations on BCD
- A-HALS
- Projected Gradient Update and the Multiplicative update

#### 3) Find $(\mathbf{W},\mathbf{H})$ numerically fast : acceleration via extrapolation

- Recall : acceleration in single variable problem
- Accelerating NMF algorithms using extrapolation

#### Convergence of the algorithms

- Application of PALM on NMF
- Convergence condition of PALM

### Does it converge?

How to show the sequence  $\{(\mathbf{W}^k, \mathbf{H}^k)\}_{k \in \mathbb{N}}$  produced by the framework converges?

Algorithm BCD framework for  $\mathcal{P}$ 

**Input:**  $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ ,  $r \in \mathbb{N}$ , an initialization  $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ ,  $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ **Output:**  $\mathbf{W}$  and  $\mathbf{H}$ 

- 1: for  $k = 1, 2, \ldots$  do
- 2: Update $[\mathbf{W}]$

• matrix-wise projected-gradient  

$$\mathbf{W} = [\mathbf{W} - t\nabla\Phi(\mathbf{W}, \mathbf{H}))]_+$$

• vector-wise A-HALS  

$$\mathbf{w}_1 = [\mathbf{w}_1 - t(\|\mathbf{h}_1\|_2^2 \mathbf{w}_1 - \mathbf{X}_1 \mathbf{h}_1^\top)]_+$$

$$\mathbf{w}_2 = [\mathbf{w}_2 - t(\|\mathbf{h}_2\|_2^2 \mathbf{w}_2 - \mathbf{X}_2 \mathbf{h}_2^\top)]_+$$

3: Update $[\mathbf{H}]$  similarly

.

4: end for

### The problem setting

$$\min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} \Phi(x, y) = f(x) + g(y) + H(x, y)$$

- f, g are extended value functions : e.g.  $f: {\rm I\!R}^n o {\rm I\!R} \cup +\infty$
- *H* is smooth (partially Lipschitz)
- No convexity will be assumed on f, g, H

For NMF :

- x, y are  $\mathbf{W}, \mathbf{H}$
- f, g are indicator functions of the non-negative constraints  $\geq 0$
- H is the data fitting term  $\|\mathbf{X} \mathbf{W}\mathbf{H}\|$

For other models

• The 2-variable case is extendable to  $n\mbox{-variable}\ \Phi(x_1,...,x_n)$  e.g. Tri-factorization, tensors

For the BCD approach

$$\begin{aligned} x^{k+1} &\in \arg\min_{x} \Phi(x, y^k) \\ y^{k+1} &\in \arg\min_{y} \Phi(x^{k+1}, y), \end{aligned}$$

the sequence  $\left\{(x^k,y^k)\right\}_{k\in\mathbb{N}}$  converges to  ${\rm crit}\Phi$  (critical point^3 of  $\Phi$ ), if

- $\Phi$  is convex and differentiable
- $\Phi(x)$  and  $\Phi(y)$  are strictly convex
  - $\Phi$  is strictly convex if one argument is fix.
  - ► The minimizer of a strictly convex function is unique, if it exists. So strict convexity ⇒ at most one global minimum
  - If fact the strict convexity is imposed for the uniququess of solution for  $\min_x \Phi(x)$  and  $\min_y \Phi(y)$

<sup>&</sup>lt;sup>3</sup>a.k.a. stationary point.

•  $\Phi$  has to be convex and differentiable

For NMF,  $\Phi$  is not convex nor differentiable because indicator function is not smooth

Φ(x) and Φ(y) are strictly convex.
 For NMF, it means W and H are full rank.
 What if no strict convexity?

### What if not strictly convex — use proximal

Proximal term relaxes the strict convexity assumption

$$\begin{aligned} x^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x, y^{k}) + \frac{c_{k}}{2} \|x - x^{k}\|^{2} \right\}, \quad c_{k} \in \mathbb{R}_{+} \\ y^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x^{k+1}, y) + \frac{d_{k}}{2} \|y - y^{k}\|^{2} \right\}, \quad d_{k} \in \mathbb{R}_{+} \end{aligned}$$

### What if not strictly convex — use proximal

Proximal term relaxes the strict convexity assumption

$$\begin{aligned} x^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x, y^{k}) + \frac{c_{k}}{2} \|x - x^{k}\|^{2} \right\}, \quad c_{k} \in \mathbb{R}_{+} \\ y^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x^{k+1}, y) + \frac{d_{k}}{2} \|y - y^{k}\|^{2} \right\}, \quad d_{k} \in \mathbb{R}_{+} \end{aligned}$$

By adding the quadratic term (with a sufficiently large  $c_k, d_k$ ), the functions  $\Phi(x, y^k) + \frac{c_k}{2} ||x - x^k||^2$  and  $\Phi(x^{k+1}, y) + \frac{d_k}{2} ||y - y^k||^2$  are strictly convex. Also,

- Fact 1. {(x<sup>k</sup>, y<sup>k</sup>)}<sub>k∈ℕ</sub> produced by such proximal regularized iteration is non-increasing in Φ.
   Direct proof by definition : Φ(x<sup>k+1</sup>, y<sup>k+1</sup>) ≤ Φ(x<sup>k+1</sup>, y<sup>k</sup>) ≤ Φ(x<sup>k</sup>, y<sup>k</sup>).
- Fact 2.  $\left\{\Phi(x^k, y^k)\right\}_{k \in \mathbb{N}}$  is bounded below by  $\inf \Phi$ .

 $\inf \Phi \le \dots \le \Phi(x^{k+1}, y^{k+1}) \le \Phi(x^k, y^k) \le \dots \le \Phi(x^1, y^1) \le \Phi(x^0, y^0)$ 

Cauchy Sequence argument : Under fact 1 + 2 IF  $\inf \Phi > -\infty$ THEN  $\{\Phi(x^k, y^k)\}_{k \in \mathbb{N}}$  converges to a real number.

Note that it only tells  $\left\{\Phi(x^k,y^k)\right\}_{k\in\mathbb{N}}$  converge, but it does not tell where it will converge!!!!

# **Theorem (Attouch10<sup>†</sup>)** If function $\Phi$ fulfill the **Kurdyka** – **Lojasiewicz property**, all bounded<sup>4</sup> sequences generated by proximal regularized iteration converge to crit $\Phi$ .

† Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Lojasiewicz inequality
 H Attouch, J Bolte, P Redont, A Soubeyran
 Mathematics of Operations Research 35 (2), 438-457,2010

- No convexity is required.
- crit⊕ is the first order stationary point (i.e. local minima), not the global optima because ncvx problems are generally NP-Hard under numerical descent schemes.
- Good for problem that local minima are almost as good as global minima — when non-convex problem becomes not scary

 $<sup>^4</sup>$  If the seugence is not bounded then it will diverge to  $\infty$ 

### Drawback of proximal regularized GS iteration

The proximal regularized iterations produces  $\{(x^k, y^k)\}_{k\in\mathbb{N}}$  via

$$x^{k+1} \in \arg\min_{x} \left\{ \Phi(x, y^k) + \frac{c_k}{2} \|x - x^k\|^2 \right\}, \quad y^{k+1} \in \arg\min_{x} \left\{ \Phi(x^{k+1}, y) + \frac{d_k}{2} \|y - y^k\|^2 \right\}.$$

which requires exact minimization of  $\Phi(x, y) = f(x) + g(y) + H(x, y)$ . For a non-convex and non-smooth  $\Phi$ , such exact minimization of may be hard/impossible.

Naming : Proximal regularized Gauss-Seidel iteration (prGSi)

How to improve prGSi : bypass such difficulty via *approximating* prGSi via proximal linearization of each subproblem — Proximal Alternating Linearized Minimization (PALM) algorithm

Why PALM : a useful framework covers many algorithms

$${\sf Problem} \quad : \quad \Phi(x,y) = f(x) + g(y) + H(x,y)$$

Note : no constraint as they are moved into f, g.

Problem : 
$$\Phi(x,y) = f(x) + g(y) + H(x,y)$$

Note : no constraint as they are moved into f, g.

#### BCD / Gauss-Seidel iteration

$$x^{k+1} \in \arg\min_x \Phi(x, y^k), \quad y^{k+1} \in \arg\min_x \Phi(x^{k+1}, y)$$

$$\mathsf{Problem} \quad : \quad \Phi(x,y) = f(x) + g(y) + H(x,y)$$

Note : no constraint as they are moved into f, g.

#### BCD / Gauss-Seidel iteration

$$x^{k+1} \in \arg\min_{x} \Phi(x, y^k), \quad y^{k+1} \in \arg\min_{x} \Phi(x^{k+1}, y)$$

Proximal regulared GS iteration (prGSi)

$$\begin{aligned} x^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x, y^{k}) + \frac{c_{k}}{2} \|x - x^{k}\|^{2} \right\} \\ y^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x^{k+1}, y) + \frac{d_{k}}{2} \|y - y^{k}\|^{2} \right\} \end{aligned}$$

$$\mathsf{Problem} \quad : \quad \Phi(x,y) = f(x) + g(y) + H(x,y)$$

Note : no constraint as they are moved into f, g.

#### BCD / Gauss-Seidel iteration

$$x^{k+1} \in \arg\min_x \Phi(x, y^k), \quad y^{k+1} \in \arg\min_x \Phi(x^{k+1}, y)$$

Proximal regulared GS iteration (prGSi)

$$\begin{aligned} x^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x, y^{k}) + \frac{c_{k}}{2} \|x - x^{k}\|^{2} \right\} \\ y^{k+1} &\in & \arg\min_{x} \left\{ \Phi(x^{k+1}, y) + \frac{d_{k}}{2} \|y - y^{k}\|^{2} \right\} \end{aligned}$$

PALM

$$\begin{aligned} x^{k+1} &\in & \arg\min_{x} \left\{ \hat{\Phi}(x, y^{k}) + \frac{c_{k}}{2} \|x - x^{k}\|^{2} \right\} \\ y^{k+1} &\in & \arg\min_{x} \left\{ \hat{\Phi}(x^{k+1}, y) + \frac{d_{k}}{2} \|y - y^{k}\|^{2} \right\} \end{aligned}$$

i.e. PALM replaces  $\Phi$  in prGSi by approximation  $\hat{\Phi}$ 

# The linear approximation $\hat{\Phi}$ in PALM

Setting : assume H is smooth, f, g not necessarily smooth<sup>5</sup> for

$$\Phi(x,y) = f(x) + g(y) + H(x,y),$$

Recall first order Taylor approximation



PALM : i.e. approximate H by the linearized H

$$\hat{\Phi}(x, y^k) = f(x) + g(y) + \left\langle x - x^k, \nabla_x H(x^k, y^k) \right\rangle$$
$$\hat{\Phi}(x^k, y) = f(x) + g(y) + \left\langle y - y^k, \nabla_y H(x^k, y^k) \right\rangle,$$

\* In convex case Taylor approximation is under-estimator so  $\approx$  becomes  $\geq$  <sup>5</sup>If f, g include the indicator function then they are non-smooth

$$\arg\min_{x} \hat{\Phi}(x, y^{k}) = \arg\min_{x} \left\{ f(x) + g(y^{k}) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\}$$

$$\underset{x}{\arg\min} \hat{\Phi}(x, y^k) = \arg\min_{x} \left\{ f(x) + g(y^k) + \left\langle x - x^k, \nabla_x H(x^k, y^k) \right\rangle \right\}$$
$$= \arg\min_{x} \left\{ f(x) + \left\langle x - x^k, \nabla_x H(x^k, y^k) \right\rangle \right\}$$

$$\begin{aligned} \arg\min_{x} \hat{\Phi}(x, y^{k}) &= \arg\min_{x} \left\{ f(x) + g(y^{k}) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ &= \arg\min_{x} \left\{ f(x) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ \arg\min_{y} \hat{\Phi}(x^{k+1}, y) &= \arg\min_{y} \left\{ f(x^{k+1}) + g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}, \end{aligned}$$

$$\begin{aligned} \arg\min_{x} \hat{\Phi}(x, y^{k}) &= \arg\min_{x} \left\{ f(x) + g(y^{k}) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ &= \arg\min_{x} \left\{ f(x) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ \arg\min_{y} \hat{\Phi}(x^{k+1}, y) &= \arg\min_{y} \left\{ f(x^{k+1}) + g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}, \\ &= \arg\min_{y} \left\{ g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}. \end{aligned}$$

Function  $\Phi(x,y) = f(x) + g(y) + H(x,y)$  has 2 variables. So alternating minimization scheme gives

$$\begin{aligned} \arg\min_{x} \hat{\Phi}(x, y^{k}) &= \arg\min_{x} \left\{ f(x) + g(y^{k}) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ &= \arg\min_{x} \left\{ f(x) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ \arg\min_{y} \hat{\Phi}(x^{k+1}, y) &= \arg\min_{y} \left\{ f(x^{k+1}) + g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}, \\ &= \arg\min_{y} \left\{ g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}. \end{aligned}$$

i.e. we have

$$\begin{aligned} \arg\min_{x} \hat{\Phi}(x, y^{k}) &= \arg\min_{x} \left\{ f(x) + \left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle \right\} \\ \arg\min_{y} \hat{\Phi}(x^{k+1}, y) &= \arg\min_{y} \left\{ g(y) + \left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle \right\}. \end{aligned}$$

### Proximal operator in PALM

Add proximal term and apply proximal operator on 
$$\Phi = \underbrace{f+g}_{x} + H$$
:  

$$x_{k} = \arg \min_{x} \left\{ f(x) + \underbrace{\left\langle x - x^{k}, \nabla_{x} H(x^{k}, y^{k}) \right\rangle + \frac{c_{k}}{2} \|x - x^{k}\|_{2}^{2}}_{\psi_{x}} \right\}$$

$$y_{k} = \arg \min_{y} \left\{ g(y) + \underbrace{\left\langle y - y^{k}, \nabla_{y} H(x^{k+1}, y^{k}) \right\rangle + \frac{d_{k}}{2} \|y - y^{k}\|_{2}^{2}}_{\psi_{y}} \right\}$$

Like FISTA, minimizer of the smooth parts  $\boldsymbol{\psi}$  is the gradient step :

set 
$$\nabla_x \psi_x = 0$$
 get  $x = x^k - \frac{1}{c_k} \nabla_x H(x^k, y^k)$ ,  $c_k > 0$   
set  $\nabla_y \psi_y = 0$  get  $y = y^k - \frac{1}{d_k} \nabla_y H(x^{k+1}, y^k)$ ,  $d_k > 0$   
From theory of gradient descent,  $(c_k, d_k)$  can be set to be the partial  
Lipschitz constant of  $\nabla H$  98 / 11

### Proximal operator in PALM

Apply proximal operator on the non-smooth parts we have

$$\begin{array}{ll} x_k &\in & \operatorname{prox}_{f,c_k} \left( x^k - \frac{1}{c_k} \nabla_x H(x^k, y^k) \right), c_k > 0, \\ y_k &\in & \operatorname{prox}_{g,d_k} \left( y^k - \frac{1}{d_k} \nabla_x H(x^{k+1}, y^k) \right), d_k > 0. \end{array}$$

Recall, at a point u, the poximal map associated to a function  $\sigma(x)$  is

$$\operatorname{prox}_{\sigma,t}(u) = \arg\min_{x} \left\{ \sigma(x) + \frac{t}{2} \|x - u\|_2^2 \right\}$$

The standard gradient descent step  $x^k - \frac{1}{c_k} \nabla_x H(x^k, y^k)$  is the "forward step". The proximal step is the "backward step".

Therefore PALM = alternating proximal forward backward method

### Short summary

Problem : minimize  $\Phi(x,y) = f(x) + g(y) + H(x,y)$ . Starts wtih  $(x^0, y^0) \in \text{dom}\Phi$ , PALM generate  $(x^k, y^k)$  as

$$\begin{aligned} x_k &\in \operatorname{prox}_{f,c_k} \left( x^k - \frac{1}{c_k} \nabla_x H(x^k, y^k) \right), \\ y_k &\in \operatorname{prox}_{g,d_k} \left( y^k - \frac{1}{d_k} \nabla_x H(x^{k+1}, y^k) \right), \end{aligned}$$

for some  $\gamma_{1,2} > 1$ , the parameters  $c_k$ ,  $d_k$  are selected as

$$c_k = \gamma_1 L_1(y^k), \quad d_k = \gamma_2 L_2(x^{k+1}).$$

In words :

- ullet on x, perform gradient update on the smooth part of  $\Phi$
- ullet on x, perform proximal update on the non-smooth part of  $\Phi$
- $\bullet\,$  on y, perform gradient update on the smooth part of  $\Phi\,$
- on y, perform proximal update on the non-smooth part of  $\Phi$
- $c_k, d_k$  are partial Lipschitz constants of H magnified

Recall f, g in  $\Phi(x, y) = f(x) + g(x) + H(x, y)$  are extended valued.

Consider the non-regularized NMF problem

$$\mathsf{NMF}: \Phi(X,Y) = \frac{1}{2} \|M - XY\|_F^2, \ X \ge 0, Y \ge 0$$

Non-negativity constraint represented by indicator function

$$\mathcal{I}_{X \ge 0}(X) = \begin{cases} 0 & X \ge 0\\ \infty & X < 0 \end{cases} \qquad \mathcal{I}_{Y \ge 0}(Y) = \begin{cases} 0 & Y \ge 0\\ \infty & Y < 0 \end{cases}$$

Recall f,g in  $\Phi(x,y) = f(x) + g(x) + H(x,y)$  are extended valued.

Consider the non-regularized NMF problem

$$\mathsf{NMF}: \Phi(X,Y) = \frac{1}{2} \|M - XY\|_F^2, \ X \ge 0, Y \ge 0$$

Non-negativity constraint represented by indicator function

$$\mathcal{I}_{X \ge 0}(X) = \begin{cases} 0 & X \ge 0\\ \infty & X < 0 \end{cases} \qquad \mathcal{I}_{Y \ge 0}(Y) = \begin{cases} 0 & Y \ge 0\\ \infty & Y < 0 \end{cases}$$

NMF problem in unconstrained form

$$\arg\min_{X,Y} \Phi(X,Y) = \frac{1}{2} \|M - XY\|_F^2 + \mathcal{I}_{X \ge 0}(X) + \mathcal{I}_{Y \ge 0}(Y)$$

#### NMF in unconstrained form

$$\arg\min_{X,Y} \Phi(X,Y) = \underbrace{\frac{1}{2} \|M - XY\|_F^2}_{H(X,Y)} + \underbrace{\mathcal{I}_{X \ge 0}(X)}_{\text{non-smooth } f} + \underbrace{\mathcal{I}_{Y \ge 0}(Y)}_{\text{non-smooth } g}$$

#### NMF in unconstrained form

$$\arg\min_{X,Y} \Phi(X,Y) = \underbrace{\frac{1}{2} \|M - XY\|_F^2}_{H(X,Y)} + \underbrace{\mathcal{I}_{X \ge 0}(X)}_{\text{non-smooth } f} + \underbrace{\mathcal{I}_{Y \ge 0}(Y)}_{\text{non-smooth } g}$$

PALM gives

$$\begin{array}{lll} X^{k+1} & \in & \operatorname{prox}_{\mathcal{I}_{X \ge 0}, c_k} \left( X^k - \frac{1}{c_k} \nabla_X H(X^k, Y^k) \right) \\ Y^{k+1} & \in & \operatorname{prox}_{\mathcal{I}_{Y \ge 0}, d_k} \left( Y^k - \frac{1}{d_k} \nabla_Y H(X^{k+1}, Y^k) \right) \end{array}$$

#### NMF in unconstrained form

$$\arg\min_{X,Y} \Phi(X,Y) = \underbrace{\frac{1}{2} \|M - XY\|_F^2}_{H(X,Y)} + \underbrace{\mathcal{I}_{X \ge 0}(X)}_{\text{non-smooth } f} + \underbrace{\mathcal{I}_{Y \ge 0}(Y)}_{\text{non-smooth } g}$$

PALM gives

$$\begin{array}{lll} X^{k+1} & \in & \operatorname{prox}_{\mathcal{I}_{X \ge 0}, c_k} \left( X^k - \frac{1}{c_k} \nabla_X H(X^k, Y^k) \right) \\ Y^{k+1} & \in & \operatorname{prox}_{\mathcal{I}_{Y \ge 0}, d_k} \left( Y^k - \frac{1}{d_k} \nabla_Y H(X^{k+1}, Y^k) \right) \end{array}$$

Fact : proximal operator of indicator function of a convex set = projection. We recovered the alternating projected gradient methods :

$$X^{k+1} = \left[ X^k - \frac{1}{c_k} \nabla_X H(X^k, Y^k) \right]_+$$
$$Y^{k+1} = \left[ Y^k - \frac{1}{d_k} \nabla_Y H(X^{k+1}, Y^k) \right]_+$$

### Convergence condition of PALM

**Theorem (Bolte14)** For  $\Phi(x,y) = f(x) + g(y) + H(x,y)$ , sequence produced by PALM converges to a stationary point of  $\Phi$  if :

#### Assumption 1

- $f: \mathbb{R}^n \to (-\infty + \infty]$  and  $g: \mathbb{R}^m \to (-\infty + \infty]$  are proper and lower semicontinuous
- $H: \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$  is a  $C^1$ /smooth function

#### Assumption 2

- $\inf_{\mathbb{R}^n \times \mathbb{R}^m} \Phi > -\infty$ ,  $\inf_{\mathbb{R}^n} f > -\infty$ ,  $\inf_{\mathbb{R}^m} g > -\infty$
- Partial gradient  $abla_x H(x,y)$  is globally Lipschitz with  $L_1(y)$  :

$$\|\nabla_x H(x_1, y) - \nabla_x H(x_2, y)\| \le L_1(y) \|x_1 - x_2\| \forall x_1, x_2 \in \mathbb{R}^n$$

• Partial gradient  $\nabla_y H(x,y)$  is globally Lipschitz with  $L_2(x)$  :

$$\|\nabla_x H(x, y_1) - \nabla_x H(x, y_2)\| \le L_2(x) \|y_1 - y_2\| \forall y_1, y_2 \in \mathbb{R}^n$$

• Lipschitz modulus  $L_1(y^k), L_2(x^k)$  are bounded

$$L_1^{\min} \le L_1(y^k) \le L_1^{\max}, \ \ L_2^{\min} \le L_2(x^k) \le L_2^{\max}, \forall k$$

• abla H is Lipschitz on bounded subsets of  ${
m I\!R}^n imes {
m I\!R}^m$ 

$$\left\| \left( \nabla_x H(x_1, y_1) - \nabla_x H(x_2, y_2), \nabla_y H(x_1, y_1) - \nabla_y H(x_2, y_2) \right) \right\| \le M \| (x_1 - x_2, y_1 - y_2) \|$$

#### Assumption 3 $\Phi$ satisfies Kurdyka-Lojasiewicz property

J Bolte, S Sabach, M Teboulle, Proximal alternating linearized minimization or nonconvex and nonsmooth problems, Mathematical Programming 146 (1-2), 459-494, 2014  $106 \ / \ 119$ 

### Convergence condition of PALM - in words

For  $\Phi(x,y)=f(x)+g(y)+H(x,y),$  sequence produced by PALM converges to a stationary point of  $\Phi$  if

- f, g are proper, lower semicontinuous, lower bounded, extended value
- *H* is smooth such that
  - All partial gradients are globally Lipschitz with L<sub>1,2</sub>
  - All Lipschitz constants  $L_1(y^k), L_2(x^k)$  are bounded
  - $\nabla H$  is Lipschitz on bounded subsets of  $\mathbb{R}^n \times \mathbb{R}^m$
- $\Phi$  is a Kurdyka-Lojasiewicz function

### Convergence condition of PALM - in words

For  $\Phi(x,y)=f(x)+g(y)+H(x,y),$  sequence produced by PALM converges to a stationary point of  $\Phi$  if

- f, g are proper, lower semicontinuous, lower bounded, extended value
- *H* is smooth such that
  - All partial gradients are globally Lipschitz with L<sub>1,2</sub>
  - All Lipschitz constants  $L_1(y^k), L_2(x^k)$  are bounded
  - $\blacktriangleright ~ \nabla H$  is Lipschitz on bounded subsets of  ${\rm I\!R}^n \times {\rm I\!R}^m$
- $\Phi$  is a Kurdyka-Lojasiewicz function

**Theorem (Bolte14)** Let  $\{z^k\}_{k\in\mathbb{N}} = \{x^k, y^k\}_{k\in\mathbb{N}}$  be the sequence generated by PALM which is assumed to be bounded, if the above are true, then (1) The path of the sequence  $\{z^k\}_{k\in\mathbb{N}}$  has finite length

$$\sum_{k=1}^{\infty} \|z^{k+1} - z^k\| < \infty.$$

(2) The sequence  $\{z^k\}_{k \in \mathbb{N}}$  converges to a stationary point  $z^*$  of  $\Phi$ .

For the details, see the paper. Next slides give the rough idea of what is going on.

J Bolte, S Sabach, M Teboulle, Proximal alternating linearized minimization or nonconvex and nonsmooth problems, Mathematical Programming 146 (1-2), 459-494, 2014 My digestion of the proof : https://angms.science/doc/NCVX/PALM1.pdf
#### Simplified setting

Consider a function  $\Phi: \mathbb{R}^n \to (-\infty, +\infty]$  is proper, l.s.c., lower bounded. For  $z = \{x, y\}$ , let the problem be

$$(P) \quad \inf \left\{ \Phi(z) : z \in \mathbb{R}^n \times \mathbb{R}^m \right\}$$

Assume there is an algorithm  ${\cal A}$  produces a sequence  $\{z^k\}$  as  $z^{k+1}\in {\cal A}(z^k)\ ,\ k\in {\rm I\!N}$ 

Goal : prove  $z^\infty$  converge to a stationary point  $z^*$  of  $\Phi.$ 

#### Simplified setting

Consider a function  $\Phi: \mathbb{R}^n \to (-\infty, +\infty]$  is proper, l.s.c., lower bounded. For  $z = \{x, y\}$ , let the problem be

$$(P) \quad \inf \left\{ \Phi(z) : z \in \mathbb{R}^n \times \mathbb{R}^m \right\}$$

Assume there is an algorithm  $\mathcal A$  produces a sequence  $\{z^k\}$  as

$$z^{k+1} \in \mathcal{A}(z^k) \ , \ k \in \mathbb{N}$$

Goal : prove  $z^\infty$  converge to a stationary point  $z^*$  of  $\Phi.$  More precisely, to prove

$$\lim_{k\to\infty}\mathsf{dist}\Big(z^k,\omega(z^0)\Big)=0\,\,,\,\,\omega(z^0)\subset\mathsf{crit}\Phi$$

Notes :

- Why "set of stationary point" but not "a single stationary point" : here problem are ncvx, can have several local minimum !
- Algorithm A has to be a proximal method : no constraint on z as constraints are in form of indicator function in Φ = f + g + H, so objective is non-smooth ⇒ has to use proximal method !110 / 119

#### In other words

#### Given

- $\Phi: {\rm I\!R}^n \to (-\infty, +\infty]$  is proper, I.s.c., lower bounded
- Problem (P) inf  $\left\{ \Phi(z) : z \in \mathbb{R}^n \right\}$
- Algorithm  ${\mathcal A}$  producing a sequence  $\{z^k\}$  as  $z^{k+1}\in {\mathcal A}(z^k),\ k\in{\rm I\!N}$

#### Goal : prove

$$\lim_{k\to\infty}{\rm dist}\big(z^k,\omega(z^0)\big)=0\ ,\ \omega(z^0)\subset{\rm crit}\Phi$$

#### Idea : show

the trajectory of  $z^1, z^2, ..., z^k, ..., z^\infty$  has finite length.

Finite length  $\implies z^\infty$  stops at somewhere, but does not tell where

**Kurdyka-Lojasiewicz** :  $z^{\infty}$  stops at stationary point of  $\Phi$ . Why KL important : no infinite circulation in trajectory  $\implies z^{\infty}$  stops at somewhere.

#### The finite path length argument



Idea : what is the length of such path?

- Sequence oscillation  $\iff$  path length  $=\infty$
- Path with finite length  $\iff z^\infty$  stops at some where
- Stop at where : ciritical point of  $\Phi(z^0)$
- Geometrically, preventing osciallation can be achieved by semi-algebaric or deformed sharp function

#### Kurdyka-Lojasiewicz (KL) function

What is Kurdyka-Lojasiewicz (KL) condition (in formal) : KL function is a class of function that can guarantee an iterative algorithm such as gradient method or near-point method does not have a circulation orbit and converges to any stationary point.

In other words, if a function  $\Phi$  meets the KL condition, it can say up to speed of convergence to a stationary point.

How to test a function satisfies the KL condition : no need. The works of Lojasiewicz and Kurdyka show that many functions fulfill the condition.

Basically standard functions used in machine learning are all KL functions. Basically you don't need to worry too much about the KL thing.

#### So you want to study Kurdyka-Lojasiewicz condition ...

Entrance level :

Proximal alternating linearized minimization or nonconvex and nonsmooth problems

J Bolte, S Sabach, M Teboulle

Mathematical Programming 146 (1-2), 459-494, 2014

Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward – backward splitting, and regularized Gauss – Seidel methods

H Attouch, J Bolte, BF Svaiter

Mathematical Programming 137 (1-2), 91-129, 2013

- Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Lojasiewicz inequality H Attouch, J Bolte, P Redont, A Soubeyran Mathematics of Operations Research 35 (2), 438-457,2010
- Inertial proximal alternating linearized minimization (iPALM) for nonconvex and nonsmooth problems

T Pock, S Sabach

SIAM Journal on Imaging Sciences, 2016

Relationships : (2,3) are the basis of (1), (4) is accelerated version of (1),

## PALM convergence theorem applied on plain NMF algorithm

Using the convergence theorem of (Bolte14) :

given starting point  $(X^0, Y^0) \in \operatorname{dom}\Phi$ , for the NMF problem

$$\arg\min_{X,Y} \Phi_{\mathsf{NMF}}(X,Y) = \underbrace{\frac{1}{2} \|M - XY\|_F^2}_{H(X,Y)} + \mathcal{I}_{X \ge 0}(X) + \mathcal{I}_{Y \ge 0}(Y)$$

where  $\Phi_{\text{NMF}}$  satisfies KL (and other assumptions), the sequence  $\{X^k, Y^k\}_{k \in \mathbb{N}}$  generated by the alternating projected gradient (PALM)

$$X^{k+1} = \left[ X^k - \frac{1}{c_k} \nabla_X H(X^k, Y^k) \right]_+$$
  
$$Y^{k+1} = \left[ Y^k - \frac{1}{d_k} \nabla_Y H(X^{k+1}, Y^k) \right]_+$$

converge to a stationary point of  $\Phi(X^0, Y^0)$ .

## PALM convergence theorem applied on plain NMF algorithm with A-HALS

Using the convergence theorem of (Bolte14) :

given starting point  $(X^0,Y^0)\in {\rm dom}\Phi,$  for the NMF problem

$$\arg\min_{X,Y} \Phi_{\mathsf{NMF}}(X,Y) = \underbrace{\frac{1}{2} \|M - XY\|_F^2}_{H(X,Y)} + \mathcal{I}_{X \ge 0}(X) + \mathcal{I}_{Y \ge 0}(Y)$$

where  $\Phi_{\text{NMF}}$  satisfies KL (and other assumptions), the sequence  $\{X^k, Y^k\}_{k \in \mathbb{N}}$  generated by the A-HALS algorithm (PALM with repeated loop of cyclic indexing) converge to a stationary point of  $\Phi(X^0, Y^0)$ .

Convergence theorem of the extrapolated NMF algorithm

# NMF with Extrapolation does not fit in the PALM framework.

### Need other tools.

No convergence theorem so far :(

Problem formulation

$$\Phi(x,y) = f(x) + g(y) + H(x,y)$$

PALM iterations

$$\begin{array}{lll} x_k & \in & \operatorname{prox}_{f,c_k} \left( x^k - \frac{1}{c_k} \nabla_x H(x^k,y^k) \right), \\ y_k & \in & \operatorname{prox}_{g,d_k} \left( y^k - \frac{1}{d_k} \nabla_x H(x^{k+1},y^k) \right), \end{array}$$

where  $c_k = \gamma_1 L_1(y^k)$ ,  $d_k = \gamma_2 L_2(x^{k+1})$  some  $\gamma_{1,2} > 1$ .

- Condition on  $\Phi$  that sequence produced PALM converges to a stationary point.
- Examples of PALM on various applications.

- What is Non-negative Matrix Factorization, Why NMF
- How to solve NMF minimization problem
- Convergence of the NMF algorithm : PALM framework
- How to solve NMF fast with extrapolation

A.-Gillis, "Accelerating Non-negative matrix factorization by extrapolation", *Neural Computation*, Feb, 2019.

• Convergence of the NMF algorithm with extrapolation

END OF PRESENTATION.

Slide, code, preprint in angms.science

ACK : my boss Nicolas Gillis, European Research Council Grant #679515.