# Accelerating Nonnegative-X by extrapolation
### where X ∈ {Least Square, Matrix Factorization, Tensor Factorization}

Andersen Ang

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

April 1, 2019

Department of Electrical and Electronic Engineering
University of Hong Kong

# Quel est le sujet de cette présentation

- Algorithmes numérique pour les problèmes d'optimisation non linéaire avec contraintes non-négativité

- Algorithmes de Factorisation Non-Négative de Matrices/Tenseurs

# Niet-negatieve matrixfactorisatie

De algemene formulering van een matrixfactorisatie is:

$$\mathbf{X} = \mathbf{WH},$$

waarbij de dimensies als volgt zijn: $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{W} \in \mathbb{R}^{m \times r}$, en $\mathbf{H} \in \mathbb{R}^{r \times n}$.

Voor de op voorhand zelf te specificeren dimensie $r$ geldt:

$$0 < r < \min\{n, m\}.$$

# Happy April Fool's day !!!!!!!!!!!!

# Outline

## Non-negative Matrix Factorization (NMF)

**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 $\mathbf{X} = \mathbf{W}\mathbf{H}$.
- Everything is <span style="color:red">non-negative</span>.



Notation : we use $\mathbf{WH}$ instead of $\mathbf{WH}^\top$.

Exact NMF: given $(\mathbf{X}, r)$, find $(\mathbf{W}, \mathbf{H})$ s.t. $\mathbf{X} = \mathbf{W}\mathbf{H}$.

It is NP-hard (Vavasis, 2007).

Vavasis, "On the complexity of nonnegative matrix factorization", SIAM J. Optim.
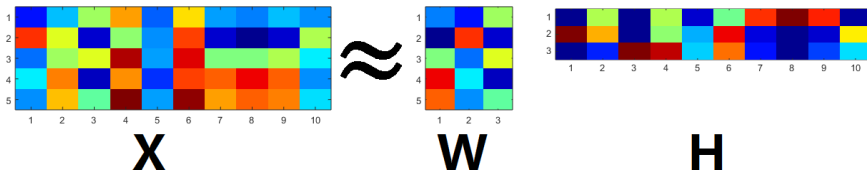
# Exact and approximate NMF

**Exact NMF**: given $(\mathbf{X}, r)$, find $(\mathbf{W}, \mathbf{H})$ s.t. $\mathbf{X} = \mathbf{WH}$.

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

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

This talk : (Low-rank) *approximate NMF*

$$\mathbf{X} \approx \mathbf{WH}, \ 1 \leq r \leq \min\{m, n\}.$$



**X** ≈ **W** **H**

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

We solve

$$[\mathbf{W}, \ \mathbf{H}] = \underset{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F.$$

- Minimizing the distance between $\mathbf{X}$ and the approximator $\mathbf{W}\mathbf{H}$ in F-norm[†].
- $\geq$ is element-wise (not positive semi-definite).
- Such minimization problem is
  - Bi-variate : two variables
  - Non-convex but block-convex (strongly convex/strictly convex)
  - Non-smooth : on the boundary between $\mathbb{R}_+$ and $\mathbb{R}_-$
  - Ill-posed and NP-hard (Vavasis, 2007)

---

[†]This talk does not consider other distance functions.

Given $(\mathbf{x}_j \in \mathbb{R}_+^m, \mathbf{W} \in \mathbb{R}_+^{m \times r})$, find $\mathbf{h}_j \in \mathbb{R}_+^r$ s.t. $\mathbf{x}_j \approx \mathbf{W}\mathbf{h}_j$ via solving

$$\mathbf{h}_j = \operatorname*{argmin}_{\mathbf{h} \geq \mathbf{0}} \|\mathbf{x}_j - \mathbf{W}\mathbf{h}_j\|_2.$$

- $\geq$ is element-wise.
- Such minimization problem is
  - ▸ Single variable
  - ▸ Non-smooth : on the boundary between $\mathbb{R}_+$ and $\mathbb{R}_-$
  - ▸ No analytic solution
  - ▸ Convex : depends on $\mathbf{W} \rightarrow$ strongly convex / strictly convex
  - ▸ Global minimizer "obtainable"
- In more "standard" notation

$$(\text{NNLS}) \quad \mathbf{x} = \operatorname*{argmin}_{\mathbf{x} \geq \mathbf{0}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$$

Given $(\mathbf{X} \in \mathbb{R}_+^{I \times J \times K}, r \in \mathbb{N})$
Find $\mathbf{U} \in \mathbb{R}_+^{I \times r}, \mathbf{V} \in \mathbb{R}_+^{J \times r}$ and $\mathbf{W} \in \mathbb{R}_+^{K \times r}$ s.t.
$\mathbf{X} \approx \mathbf{U} * \mathbf{V} * \mathbf{W}$ via solving

$$\mathbf{X} = \operatorname*{argmin}_{\{\mathbf{U}, \mathbf{V}, \mathbf{W}\} \geq \mathbf{0}} \|\mathbf{X} - \mathbf{U} * \mathbf{V} * \mathbf{W}\|_F.$$

- $\geq$ is element-wise.
- Such minimization problem is
  - Three variables
  - Non-smooth : on the boundary between $\mathbb{R}_+$ and $\mathbb{R}_-$
  - Non-convex but block-convex (strongly convex/strictly convex)
  - Global minimizer "obtainable" – unique solution under some conditions on $I, J, K, r$

---

†This talk does not consider other tensor norm.

## The scope of this talk : computation of NNLS, NMF, NTF

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

**Non-keywords** : Sparsity, Regularization, Applications, Non-negative rank, Extended Formulations, Separability, NP-Hardness

**Focus** : single-machine, serial, deterministic algorithm

**Non-focus** : multi-machine, parallel, distributed, stochastic algorithm

# 5 slides on why NMF

- ## Interpretability

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

  on non-negative data.

- ## Model correctness
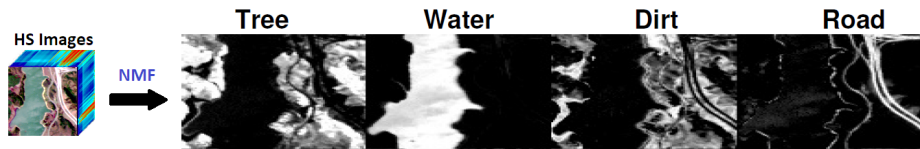
  NMF can find ground truth (under certain conditions).

- ## Mathematical curiosity

  NMF is related to some serious problems in mathematics.

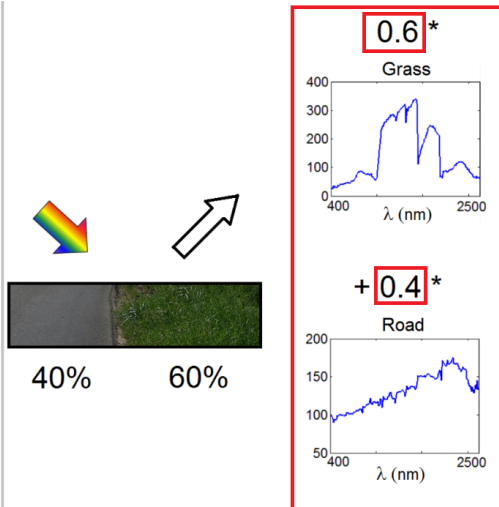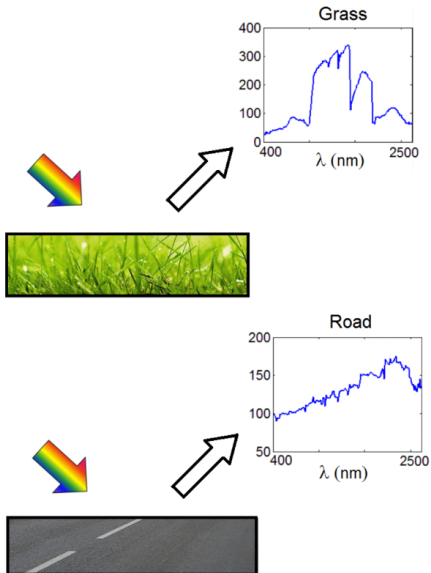- ## ~~My boss tell me to do it.~~

NMF gives good *unsupervised* image segmentation[1]



Hyper-spectral image decomposition. Figure from Zhu, F. et al., "Spectral unmixing via data-guided sparsity." IEEE Trans. Image Processing, 2014
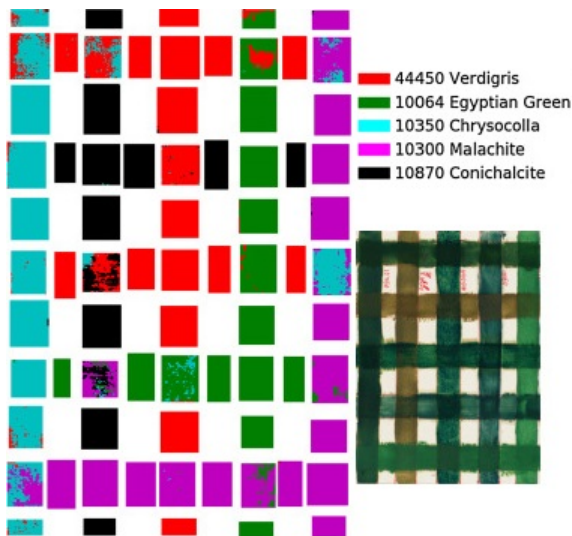
# Comment est-ce possible ?!

---

[1]Modern fancy name : "super resolution"

Hyper-spectral imaging. Figure modified from the slide of Nicolas Gillis.

# Why NMF - art work preservation example



Art work preservation. Figure from Grabowski, Bartosz, et al. "Automatic pigment identification from hyperspectral data." J. Cultural Heritage 31 (2018): 1-12.

# 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
- Audio
- 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 $\text{rank}^+ :=$ smallest $r$ such that

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

  How to find / estimate / bound $\text{rank}^+$, e.g. $\text{rank}_{\text{psd}}(\mathbf{X}) \leq \text{rank}^+(\mathbf{X})$.
- Extended formulations and combinatorics
- Log-rank Conjecture of communication system
- 3-SAT, Exponential time hypothesis, $\mathbf{P} \neq \mathbf{NP}$

# Outline

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

$$[\mathbf{W}, \ \mathbf{H}] = \underset{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}}{\mathrm{argmin}} \ \Phi(\mathbf{W}, \mathbf{H}) = \|\mathbf{X} - \mathbf{WH}\|_F.$$

- Equivalent objective function : $\frac{1}{2}\|\mathbf{X} - \mathbf{WH}\|_F^2$.

- Simplify notation : hide some $\geq \mathbf{0}, \frac{1}{2}, F$ and write

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

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

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

Approach : BCD (Block Coordinate Descent)[2]

---

**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[$\mathbf{W}$] : do something with $\Phi, \mathbf{X}, \mathbf{W}, \mathbf{H}$.
3:   Update[$\mathbf{H}$] : do something with $\Phi, \mathbf{X}, \mathbf{W}, \mathbf{H}$.
4: **end for**

---

[2]Other names : Gauss-Seidel iteration, alternating minimization (for $2$ blocks)

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

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

Approach : BCD (Block Coordinate Descent)[2]

---

**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[$\mathbf{W}$] : do something with $\Phi, \mathbf{X}, \mathbf{W}, \mathbf{H}$.
3:     Update[$\mathbf{H}$] : do something with $\Phi, \mathbf{X}, \mathbf{W}, \mathbf{H}$.
4: **end for**

---

The goal of "do something" is to achieve

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

(Actually non-increasing condition is not enough, need sufficient decrease condition)

---

[2]Other names : Gauss-Seidel iteration, alternating minimization (for $2$ blocks)

## Example 1 : alternating minimization

**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[$\mathbf{W}$] as $\mathbf{W} \leftarrow \underset{\mathbf{W} \geq 0}{\operatorname{argmin}} \Phi(\mathbf{W}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$.
3:     Update[$\mathbf{H}$] as $\mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \Phi(\mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$.
4: **end for**

## Example 1 : alternating minimization

**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[$\mathbf{W}$] as $\mathbf{W} \leftarrow \underset{\mathbf{W} \geq 0}{\mathrm{argmin}} \Phi(\mathbf{W}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$.
3:    Update[$\mathbf{H}$] as $\mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\mathrm{argmin}} \Phi(\mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$.
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 $\mathbf{W}$ (or $\mathbf{H}$) : we
have to handle them separately.

## Example 2: alternating gradient descent

---

**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** $k = 1, 2, \dots$ **do**
2:    Update[$\mathbf{W}$] as

$$\mathbf{W} \leftarrow \underset{\mathbf{W} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{U}\mathbf{H}\|_F^2 + \langle \mathbf{W} - \mathbf{U}, \nabla \Phi(\mathbf{U}) \rangle + \frac{1}{2t}\|\mathbf{U} - \mathbf{W}\|^2.$$

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

$$\mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{V}\|_F^2 + \langle \mathbf{H} - \mathbf{V}, \nabla \Phi(\mathbf{V}) \rangle + \frac{1}{2t}\|\mathbf{V} - \mathbf{H}\|^2.$$

4: **end for**

---

**Local quadratic model** : gradient descent minimizes the local quadratic model of the original objective function

$$\text{Update}[\mathbf{H}] : \ \mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$$

1. Block partitions : on how coordinate is being defined[†].
   Now : coordinate is $\mathbf{H}$ (matrix) or $\mathbf{H}(i,:)$ (vector).

$$\text{Update}[\mathbf{H}] : \ \mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

1. Block partitions : on how coordinate is being defined[†].
   Now : coordinate is $\mathbf{H}$ (matrix) or $\mathbf{H}(i, :)$ (vector).

2. Index selection (indexing) : on how coordinate is being selected[#].
   Now : cyclic indexing and A-HALS.

# Variations on BCD

$$\text{Update}[\mathbf{H}] : \ \mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$$

1. Block partitions : on how coordinate is being defined[†].
   Now : coordinate is $\mathbf{H}$ (matrix) or $\mathbf{H}(i, :)$ (vector).

2. Index selection (indexing) : on how coordinate is being selected[#].
   Now : cyclic indexing and A-HALS.

3. Update scheme : on how coordinate is being updated[#].
   Now : "exact" coordinate minimization using 1st order method (e.g. gradient descent).
   Exact = working on the original original objective function, no modification.

   Inexact = working on modified objective function. e.g. consider relaxation.

$$\text{Update}[\mathbf{H}]: \ \mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$$

1. Block partitions : on how coordinate is being defined[†].
   Now : coordinate is $\mathbf{H}$ (matrix) or $\mathbf{H}(i,:)$ (vector).

2. Index selection (indexing) : on how coordinate is being selected[#].
   Now : cyclic indexing and A-HALS.

3. Update scheme : on how coordinate is being updated[#].
   Now : "exact" coordinate minimization using 1st order method (e.g. gradient descent).
   Exact = working on the original original objective function, no modification.

   Inexact = working on modified objective function. e.g. consider relaxation.

4. Other variants

   † 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 $\mathbf{W}$ and row of $\mathbf{H}$), we have

$$\Phi(\mathbf{w}_i, \mathbf{h}_i) = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr}\langle \mathbf{X}_i, \mathbf{w}_i \mathbf{h}_i \rangle + c.$$

Says coordinates are vectors (col. of $\mathbf{W}$ and row of $\mathbf{H}$), we have

$$\Phi(\mathbf{w}_i, \mathbf{h}_i) = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr}\langle \mathbf{X}_i, \mathbf{w}_i\mathbf{h}_i \rangle + c.$$

**Alternating minimization using cyclic indexing**
Other name : BCD with $r = 2$ with cyclic component selection
Domain name in NMF : HALS (Hierarchical alternating least squares[†])

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

# The idea of HALS and A-HALS

Says coordinates are vectors (col. of $\mathbf{W}$ and row of $\mathbf{H}$), we have

$$\Phi(\mathbf{w}_i, \mathbf{h}_i) = \|\mathbf{w}_i\|_2^2 \|\mathbf{h}_i\|_2^2 - 2\mathrm{tr}\langle \mathbf{X}_i, \mathbf{w}_i \mathbf{h}_i \rangle + c.$$

**Alternating minimization using cyclic indexing**
Other name : BCD with $r = 2$ with cyclic component selection
Domain name in NMF : HALS (Hierarchical alternating least squares[†])

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**[#] (Accelerated-HALS)
A special kinds of cyclic coordinate selection

Update order : $\underbrace{\mathbf{w}_1 \rightarrow \mathbf{w}_2 \rightarrow \cdots \rightarrow \mathbf{w}_r}_{\text{several times!!}} \rightarrow \underbrace{\mathbf{h}_1 \rightarrow \mathbf{h}_2 \rightarrow \cdots \rightarrow \mathbf{h}_r}_{\text{several times!!}} \rightarrow ...$

† 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.

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

Projected[†] gradient descent with step size $t \geq 0$

$$\mathbf{w}_i = \mathbf{w}_i - t \underbrace{(\|\mathbf{h}_i\|_2^2 \mathbf{w}_i - \mathbf{X}_i \mathbf{h}_i^\top)}_{\nabla_{\mathbf{w}_i} \Phi}, \quad \mathbf{h}_i = \mathbf{h}_i - t \underbrace{(\|\mathbf{w}_i\|_2^2 \mathbf{h}_i - \mathbf{w}_i^\top \mathbf{X})}_{\nabla_{\mathbf{h}_i} \Phi}.$$

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

Projected[†] gradient descent with step size $t \geq 0$

$$\mathbf{w}_i = \mathbf{w}_i - t \underbrace{(\|\mathbf{h}_i\|_2^2 \mathbf{w}_i - \mathbf{X}_i \mathbf{h}_i^\top)}_{\nabla_{\mathbf{w}_i} \Phi}, \quad \mathbf{h}_i = \mathbf{h}_i - t \underbrace{(\|\mathbf{w}_i\|_2^2 \mathbf{h}_i - \mathbf{w}_i^\top \mathbf{X})}_{\nabla_{\mathbf{h}_i} \Phi}.$$

| **Algorithm** HALS | **Algorithm** A-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: $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_3^\top)$ | 5: Compute $\mathbf{C} = \mathbf{W}^\top \mathbf{W}$, $\mathbf{D} = \mathbf{W}^\top \mathbf{X}$ |
| 6: $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_3^\top \mathbf{X}_3)$ | 6: $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ |
| 7: ... | 7: $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ |
| | 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[†] gradient descent with step size $t \geq 0$

$$\mathbf{w}_i = \mathbf{w}_i - t \underbrace{(\|\mathbf{h}_i\|_2^2 \mathbf{w}_i - \mathbf{X}_i \mathbf{h}_i^\top)}_{\nabla_{\mathbf{w}_i} \Phi}, \quad \mathbf{h}_i = \mathbf{h}_i - t \underbrace{(\|\mathbf{w}_i\|_2^2 \mathbf{h}_i - \mathbf{w}_i^\top \mathbf{X})}_{\nabla_{\mathbf{h}_i} \Phi}.$$

| **Algorithm** HALS | **Algorithm** A-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 $\mathbf{A} = \mathbf{HH}^\top$, $\mathbf{B} = \mathbf{XH}^\top$ |
| 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: $\mathbf{w}_3 = \mathbf{w}_3 - t(\|\mathbf{h}_3\|_2^2 \mathbf{w}_3 - \mathbf{X}_3 \mathbf{h}_3^\top)$ | 5: Compute $\mathbf{C} = \mathbf{W}^\top \mathbf{W}$, $\mathbf{D} = \mathbf{W}^\top \mathbf{X}$ |
| 6: $\mathbf{h}_3 = \mathbf{h}_3 - t(\|\mathbf{w}_3\|_2^2 \mathbf{h}_3 - \mathbf{w}_3^\top \mathbf{X}_3)$ | 6: $\mathbf{h}_1 = \mathbf{h}_1 - t(\|\mathbf{w}_1\|_2^2 \mathbf{h}_1 - \mathbf{w}_1^\top \mathbf{X}_1)$ |
| 7: ... | 7: $\mathbf{h}_2 = \mathbf{h}_2 - t(\|\mathbf{w}_2\|_2^2 \mathbf{h}_2 - \mathbf{w}_2^\top \mathbf{X}_2)$ |
| | 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* :

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

*pre-computing* and *re-use* of these terms gain extra efficiency, improvement is significant for big data[#]" — always A-HALS!

---

†Projection step not shown here. # Even more significant in terms of BLAS if the matrices are *sparse*.

# The projected gradient descent update

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

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

How to pick the step-size ?

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

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

How to pick the step-size ?

A simple scheme $t = \dfrac{1}{L_{\Phi_{\mathbf{W}}}}$, 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 $\mathbf{W}$ :

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

How to pick the step-size ?

A simple scheme $t = \dfrac{1}{L_{\Phi_{\mathbf{w}}}}$, where $L_{\Phi_{\mathbf{w}}} = $ the Lipschitz constant of $\nabla_{\mathbf{W}}\Phi$ (smoothness constant).

- $L_{\Phi_{\mathbf{w}}} = $ largest singular value of $\mathbf{H}\mathbf{H}^\top$
- $\text{Proj}_{\mathbb{R}_+}$ is basically $[\ \cdot\ ]_+ = \max\{\cdot, 0\}$.

Hence in close form :

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

PGD update is much faster than the *Multiplicative Update*.

# Multiplicative Update

MU :

- It takes a small step size $t$ such that $\mathbf{W}^{k+1}$ stays within $\mathbb{R}_+$, no projection.

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

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

- It converges **very slowly**. In general, don't use MU.
  Why: to make sure $\mathbf{W}$ stays within $\mathbb{R}_+$, MU take small step $\implies$ slow !

PGD :

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

Relative Error vs iterations

MU = timid, shy guy that is too cautious on making mistake.
PGD = brave guy that is fine of making mistake by doing correction.
Here "mistake" = "outside $\mathbb{R}_+$", "correction" = "$\mathrm{Proj}_{\mathbb{R}_+}$".

# Outline

# Let's accelerate !

The next many slides : make PGD converges even more fast



Recall : NMF is **NP-Hard**.
What's the acceleration for : obtain a *local* solution faster

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

# Recall : acceleration in single variable problem

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

At step $k$ :

No acceleration : $\quad x_{k+1} = \mathsf{Update}[x_k]$.

With acceleration : $\quad x_{k+1} = \mathsf{Update}[y_k], \quad y_{k+1} = \mathsf{Extrapolate}[x_{k+1}, x_k]$.

# Recall : acceleration in single variable problem

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

At step $k$ :

No acceleration :   $x_{k+1} = \mathsf{Update}[x_k]$.

With acceleration :   $x_{k+1} = \mathsf{Update}[y_k]$,   $y_{k+1} = \mathsf{Extrapolate}[x_{k+1}, x_k]$.

To be specific :

$$\text{PGD Update} \quad x_{k+1} = \mathrm{Proj}_{\mathcal{C}}(x_k - t_k \nabla f(x_k)).$$
$$\text{Linear extrapolation} \quad x_{k+1} = \mathrm{Proj}_{\mathcal{C}}(y_k - t_k \nabla f(y_k)).$$
$$y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$$

i.e. $\mathsf{Extrapolate}[x_{k+1}, x_k]$ 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.



**Observations**

- When the level set of $f(\mathbf{x})$ is circular, GD goes to $\mathbf{x}^*$ very fast.
  (In fact, in 1 step GD goes to $\mathbf{x}^*$)
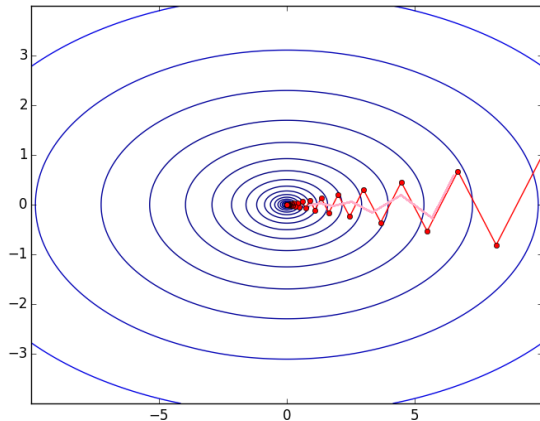- When the level set of $f(\mathbf{x})$ is elliptic, GD zigzags (and slow).

**Questions**

- Why? Where does this zigzag come from?
- How to deal with it : how to improve GD?

26 / 144

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

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



Length of pink segment $<$ length of the corresponding red segment $\implies$ points on pink segment is closer to axis $y = 0$ , gradient stronger $x$-component $\implies$ less oscillation along $y$-direction.

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

# What optimization people do to counter zig-zag?

**Do tricks on direction** : by extrapolation with momentum.



Idea : apply extrapolation.
Extrapolate = add gradient history.

(1) if gradients in consecutive steps have consistent direction
$\implies$ extrapolate = accelerate.
(2) if gradients in consecutive steps oscillates (continuously changing direction)
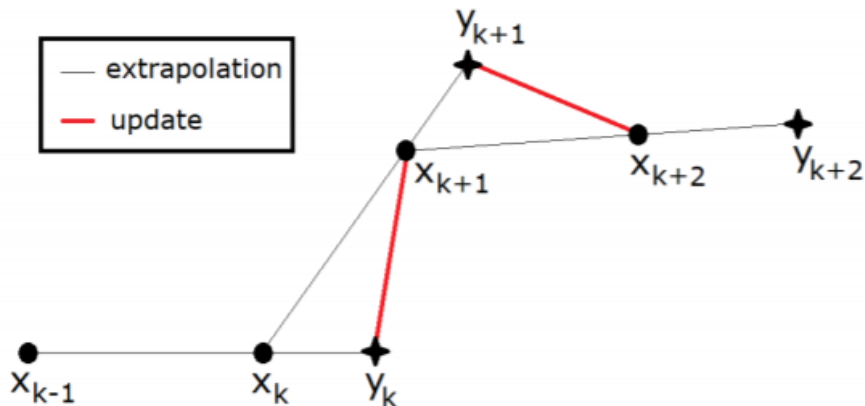$\implies$ extrapolate = damp oscillation = acceleration.

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

$$x_{k+1} = \mathsf{Update}[y_k], \quad y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$$
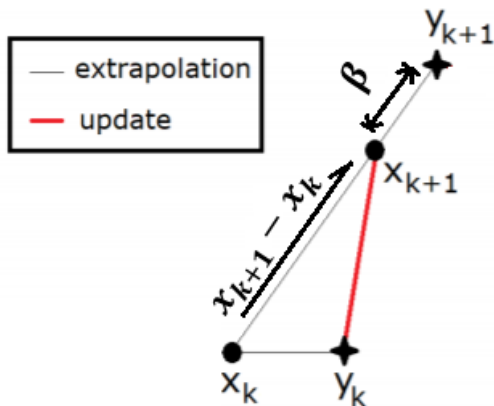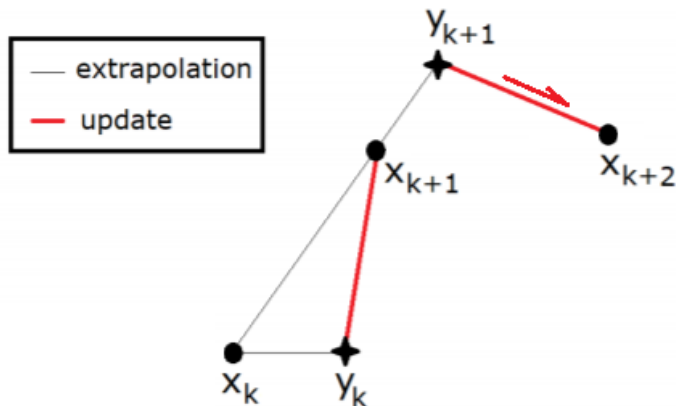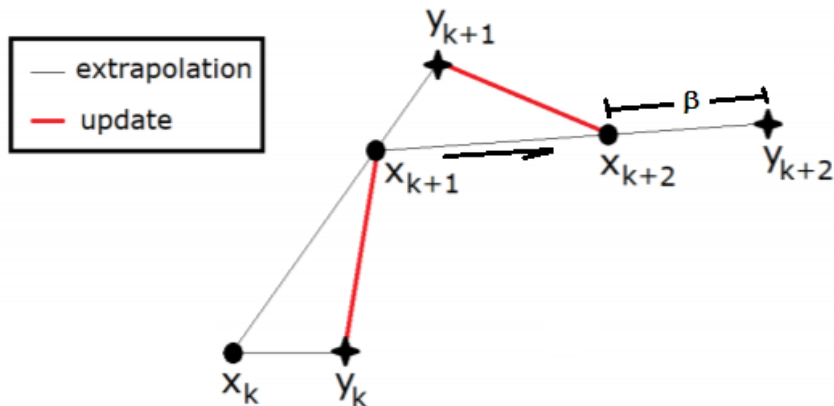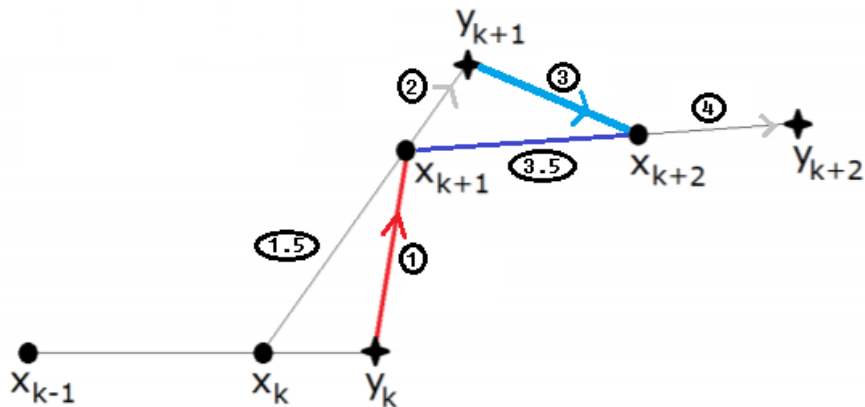
# The geometry of the extrapolation

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

# The geometry of the extrapolation

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

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

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

$$x_{k+1} = \mathsf{Update}[y_k], \quad y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$$
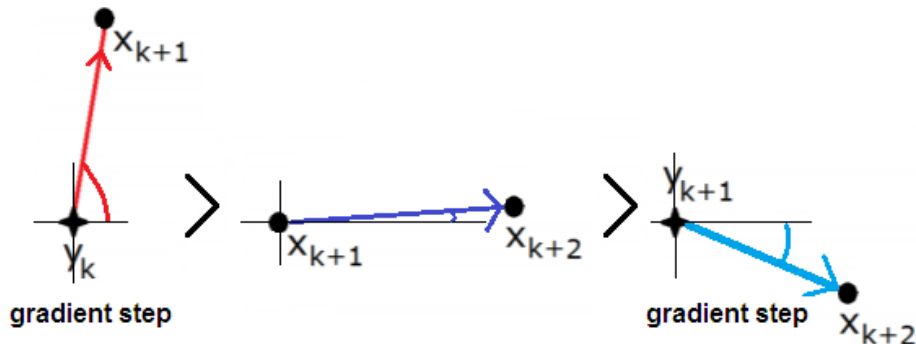
# The geometry of extrapolation

We always have

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

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



gradient step                    gradient step

# Nesterov's acceleration

For **convex** (smooth strongly-convex) function

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

$$\beta_k = \beta = \frac{1 - \sqrt{\kappa'}}{1 + \sqrt{\kappa'}}, \ \kappa' = \frac{1}{\kappa}, \ \kappa = \frac{\sigma_{\max}(\mathbf{Q})}{\sigma_{\min}(\mathbf{Q})} = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$$
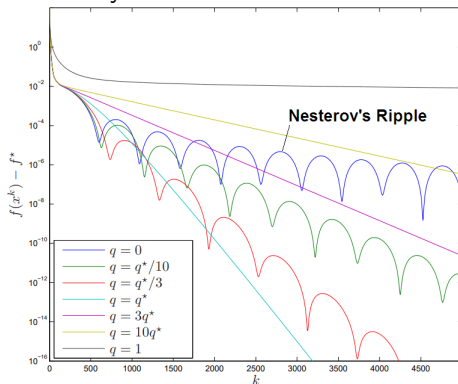
70 / 144

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 $\implies$ extrapolation may $\uparrow$objective value
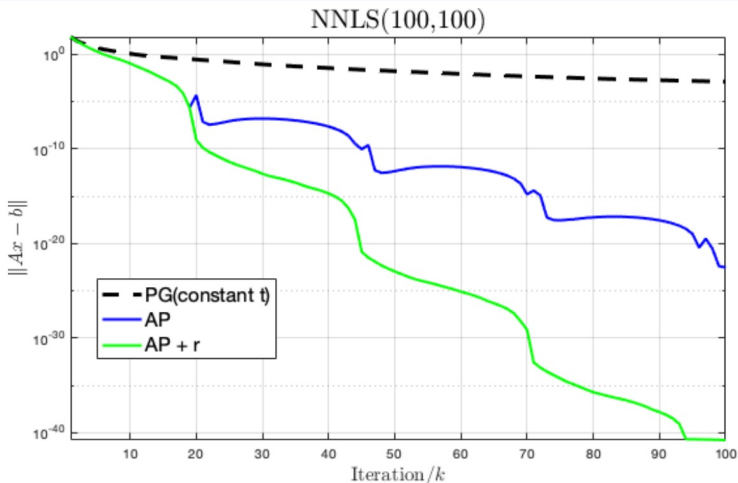- Extrapolation = a risky move



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"

NNLS(100,100)

78 / 144

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

## Our case : NMF is not cvx

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

- No strong cvx parameter, cannot use expression likes $\beta_k = \dfrac{1 - \sqrt{\kappa}}{1 + \sqrt{\kappa}}$.

- Direct application of Nesterov's $\beta$ sequence on PGD/A-HALS will give erratic convergence behaviour

  Mitchell, Drew, Nan Ye, and Hans De Sterck. "Nesterov Acceleration of Alternating Least Squares for Canonical Tensor Decomposition." arXiv:1810.05846 (2018)

# Our case : NMF is not cvx

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

- No strong cvx parameter, cannot use expression likes $\beta_k = \dfrac{1 - \sqrt{\kappa}}{1 + \sqrt{\kappa}}$.

- Direct application of Nesterov's $\beta$ sequence on PGD/A-HALS will give erratic convergence behaviour

  Mitchell, Drew, Nan Ye, and Hans De Sterck. "Nesterov Acceleration of Alternating Least Squares for Canonical Tensor Decomposition." arXiv:1810.05846 (2018)

For the acceleration scheme of the two variables

$$
\begin{cases}
\text{On } \mathbf{W} \begin{cases}
\text{Update} \quad \mathbf{W}_{\text{new}} = \text{Update}[\mathbf{Y}_{\text{old}}, \mathbf{H}_{\text{old}}] \\
\\
\text{Extrapolate} \quad \mathbf{Y}_{\text{new}} = \mathbf{W}_{\text{new}} + \beta_k^{\mathbf{W}}(\mathbf{W}_{\text{new}} - \mathbf{W}_{\text{old}})
\end{cases} \\
\text{On } \mathbf{H} \begin{cases}
\text{Update} \quad \mathbf{H}_{\text{new}} = \text{Update}[\mathbf{W}_{\text{new}}, \mathbf{G}_{\text{old}}] \\
\\
\text{Extrapolate} \quad \mathbf{G}_{\text{new}} = \mathbf{H}_{\text{new}} + \beta_k^{\mathbf{H}}(\mathbf{H}_{\text{new}} - \mathbf{H}_{\text{old}})
\end{cases}
\end{cases}
$$

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

## Approach : an ad hoc heurisitic in the "line search" style.

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

## Details of the extrapolation

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)
- If $\beta \in (0, 1)$ : extrapolation, doing risky step
- If $\beta = \{1, 0\}$ : doing {very risky, no} extrapolation
- Can't use line search[†] to find $\beta$ : experimentally found $\beta$ close to 0

  – effectively doing no extrapolation, waste resource on line search

## Details of the extrapolation

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
- Can't use line search[†] to find $\beta$ : experimentally found $\beta$ close to 0

  – effectively doing no extrapolation, waste resource on line search

In the "walking person metaphor" :

| | |
|---:|:---|
| MU | shy guy walking in caution with small step size |
| PGD | brave guy walking with reasonably step size |
| E-PGD | ambious guy walking with big step size |
| E-A-HALS | crazy guy walking with big step size in coordinate manner |

†Line search to minimize the objective function directly – performed **before** the update

# Details : Update[$\beta_k$]

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

---

**Algorithm** A dynamic **line search style**[†] **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**

---

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

[†]Line search **after** updates of $\mathbf{W}$ and $\mathbf{H}$ – performed **after** the update!

# The logic of tuning $\beta$

The idea is to update $\beta_k$ based on the increase or decrease of the objective function. Let $e^k = \Phi(\mathbf{W}^k, \mathbf{H}^k)$, then

$$
\beta_{k+1} = \begin{cases} \min\{\gamma\beta_k, \bar{\beta}\} & \text{if } e^k \leq e^{k-1} \\ \dfrac{\beta_k}{\eta} & \text{if } e^k > e^{k-1} \end{cases} \tag{1}
$$

where $\gamma > 1$, and $\eta > 1$ are constants and $\bar{\beta}_0 = 1$ with the update

# The logic of tuning $\beta$

The idea is to update $\beta_k$ based on the increase or decrease of the objective function. Let $e^k = \Phi(\mathbf{W}^k, \mathbf{H}^k)$, then

$$\beta_{k+1} = \begin{cases} \min\{\gamma\beta_k, \bar{\beta}\} & \text{if } e^k \leq e^{k-1} \\ \dfrac{\beta_k}{\eta} & \text{if } e^k > e^{k-1} \end{cases} \tag{1}$$

where $\gamma > 1$, and $\eta > 1$ are constants and $\bar{\beta}_0 = 1$ with the update

$$\bar{\beta}_{k+1} = \begin{cases} \min\{\bar{\gamma}\bar{\beta}_k, 1\} & \text{if } e^k \leq e^{k-1} \text{ and } \bar{\beta}_k < 1 \\ \beta_k & \text{if } e^k > e^{k-1} \end{cases}. \tag{2}$$

**Case 1. The error decreases** : $e^k \le e^{k-1}$

- It means the current $\beta$ value is "good"

**Case 1. The error decreases** : $e^k \leq e^{k-1}$

- It means the current $\beta$ value is "good"
- We can be more ambitious on the extrapolation
    - ▸ i.e., we increase the value of $\beta$
    - ▸ How : multiplying it with a growth factor $\gamma > 1$

$$\beta_{k+1} = \beta_k \gamma$$

**Case 1. The error decreases** : $e^k \le e^{k-1}$

- It means the current $\beta$ value is "good"
- We can be more ambitious on the extrapolation
    - ▸ i.e., we increase the value of $\beta$
    - ▸ How : multiplying it with a growth factor $\gamma > 1$

$$\beta_{k+1} = \beta_k \gamma$$

- Note that the growth of $\beta$ cannot be indefinite
    - ▸ i.e., we put a ceiling parameter $\bar{\beta}$ to upper bound the growth
    - ▸ How : use $\min$

$$\beta_{k+1} = \min\{\beta_k \gamma, \bar{\beta}_k\}$$

    - ▸ $\bar{\beta}$ itself is also updated dynamically with a growth factor $\bar{\gamma}$ with the upper bound 1.

**Case 2. The error increases** : $e^k > e^{k-1}$

- It means the current $\beta$ value is "bad" (too large)

**Case 2. The error increases** : $e^k > e^{k-1}$

- It means the current $\beta$ value is "bad" (too large)
- We become less ambitious on the extrapolation
  - i.e., we decrease the value of $\beta$
  - How : dividing it with the decay factor $\eta > 1$

$$\beta_{k+1} = \frac{\beta_k}{\eta}$$

**Case 2. The error increases** : $e^k > e^{k-1}$

- It means the current $\beta$ value is "bad" (too large)
- We become less ambitious on the extrapolation
    - i.e., we decrease the value of $\beta$
    - How : dividing it with the decay factor $\eta > 1$

$$\beta_{k+1} = \frac{\beta_k}{\eta}$$

- As $f$ is often a continuous and smooth, for $\beta_k$ being too large, such value of $\beta$ will also be too large at iteration $k+1$
    - i.e., we have to avoid $\beta_{k+1}$ to grow back to $\beta_k$ (the "bad" value) too soon
    - How : we set the ceiling parameter

$$\bar{\beta}_{k+1} = \beta_k$$

# The full algo of Accelerated NMF using extrapolation

**Input:** $\mathbf{X}$, initialization $\mathbf{W}, \mathbf{H}$, parameters $hp \in \{1, 2, 3\}$ (extrapolation/projection of $\mathbf{H}$).

**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** $k = 1, 2, \ldots$ **do**
3:    Compute $\mathbf{H}_n$ by $\min\limits_{\mathbf{H}_n \geq 0} ||\mathbf{X} - \mathbf{W}_y \mathbf{H}_n||_F^2$ using $\mathbf{H}_y$ as initial iterate.
4:      **if** $hp \geq 2$ **then**
5:          Extrapolate: $\mathbf{H}_y = \mathbf{H}_n + \beta_k(\mathbf{H}_n - \mathbf{H})$.
6:      **end if**
7:      **if** $hp = 3$ **then**
8:          Project: $\mathbf{H}_y = \max(0, \mathbf{H}_y)$.
9:      **end if**
10:     Compute $\mathbf{W}_n$ by $\min\limits_{\mathbf{W}_n \geq 0} ||\mathbf{X} - \mathbf{W}_n \mathbf{H}_y||_F^2$ using $\mathbf{W}_y$ as initial iterate.
11:     Extrapolate: $\mathbf{W}_y = \mathbf{W}_n + \beta_k(\mathbf{W}_n - \mathbf{W})$.
12:     **if** $hp = 1$ **then**
13:         Extrapolate: $\mathbf{H}_y = \mathbf{H}_n + \beta_k(\mathbf{H}_n - \mathbf{H})$.
14:     **end if**
15:     Compute error: $e(k) = ||\mathbf{X} - \mathbf{W}_n \mathbf{H}_y||_F$.
16:     **if** $e(k) > e(k-1)$ **then**
17:         Restart: $\mathbf{H}_y = \mathbf{H}_n; \mathbf{W}_y = \mathbf{W}_n$.
18:     **else**
19:         $\mathbf{H} = \mathbf{H}_n; \mathbf{W} = \mathbf{W}_n$.
20:     **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** $k = 1, 2, \ldots$ **do**
3:     **Up**date[$\mathbf{H}_n$] w.r.t. $\mathbf{H}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$ using $\mathbf{H}_y$ as initial iterate.
4:     **Up**date[$\mathbf{W}_n$] wr.t. $\mathbf{W}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$ using $\mathbf{W}_y$ as initial iterate.
5:     **Ex**trapolate[$\mathbf{W}_y$] : $\mathbf{W}_y = \mathbf{W}_n + \beta_k(\mathbf{W}_n - \mathbf{W})$.
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** $k = 1, 2, \ldots$ **do**
3:    **Up**date$[\mathbf{H}_n]$ w.r.t. $\mathbf{H}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$ using $\mathbf{H}_y$ as initial iterate.
4:    **Ex**trapolate$[\mathbf{H}_y]$ : $\mathbf{H}_y = \mathbf{H}_n + \beta_k(\mathbf{H}_n - \mathbf{H})$.
5:    **Up**date$[\mathbf{W}_n]$ wr.t. $\mathbf{W}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$ using $\mathbf{W}_y$ as initial iterate.
6:    **Ex**trapolate$[\mathbf{W}_y]$ : $\mathbf{W}_y = \mathbf{W}_n + \beta_k(\mathbf{W}_n - \mathbf{W})$.
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** $k = 1, 2, \ldots$ **do**
3:     **Up**date[$\mathbf{H}_n$] w.r.t. $\mathbf{H}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_y, \mathbf{H}_n$ using $\mathbf{H}_y$ as initial iterate.
4:     **Ex**trapolate[$\mathbf{H}_y$] : $\mathbf{H}_y = \mathbf{H}_n + \beta_k(\mathbf{H}_n - \mathbf{H})$.
5:     **Pro**ject: $\mathbf{H}_y = \max(0, \mathbf{H}_y)$.
6:     **Up**date[$\mathbf{W}_n$] wr.t. $\mathbf{W}_n \geq 0$ with $\mathbf{X}, \mathbf{W}_n, \mathbf{H}_y$ using $\mathbf{W}_y$ as initial iterate.
7:     **Ex**trapolate[$\mathbf{W}_y$] : $\mathbf{W}_y = \mathbf{W}_n + \beta_k(\mathbf{W}_n - \mathbf{W})$.

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[$\mathbf{H}_y$] | N | Extrap[$\mathbf{H}_y$] | N |
| | | | | Project[$\mathbf{H}_y$] | Y |
| Extrap[$\mathbf{H}_y$] | N | Update[$\mathbf{W}_n$] | Y | Update[$\mathbf{W}_n$] | Y |
| Extrap[$\mathbf{W}_y$] | N | Extrap[$\mathbf{W}_y$] | N | Extrap[$\mathbf{W}_y$] | N |

There are variations on the chain structure of the update, for examples

- Update $\mathbf{W}$ → extrapolate $\mathbf{W}$ → update $\mathbf{H}$ → extrapolate $\mathbf{H}$
- Update $\mathbf{W}$ → extrapolate $\mathbf{W}$ → update $\mathbf{H}$ → extrapolate $\mathbf{H}$ → project $\mathbf{H}$
- Update $\mathbf{W}$ → update $\mathbf{H}$ → extrapolate $\mathbf{W}$ → extrapolate $\mathbf{H}$

The comparisons of these three schemes : see the paper.

Iteration $k$

**Iteration $k$**

**Iteration $k$**

**Open question** : why certain structure has a better performance than others

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

(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 $\mathbf{W}_n^\infty = \mathbf{W}_y^\infty$ (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 $\sim \mathcal{U}[0,1]$, $m = n = 200$, $r = 20$
- Real $\mathbf{X}$ 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 $\log \neq 1$)

Low-rank synthetic data

Image data

Image data

Text data

Fast conclusion : E wins.

# 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[†] = an extrapolated proximal type algorithm, with convergence proof.

<div align="center">

Fast conclusion : E wins and beats APG-MF[†].

</div>

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

# Overall results : E wins!

| Method | Data | Ex wins? |
|--------|------|----------|
| A-HALS | Low/full rank synthetic data | YES |
| | **Dense Image data**[†] | YES |
| | **Sparse text data**[#] | YES |
| ANLS | Low/full rank synthetic data | YES |
| | **Dense Image data**[†] | 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 ≈ E-ANLS, although A-HALS > ANLS
  - ▸ Sparse data : E-A-HALS ≫ 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!

# Outline

## Tensor extension

(Joint-work with Jeremy E. Cohen of IRISA, Rennes, France)

Extend the idea of extrapolation to the tensor cases; more precisely to the **Non-negative Canonical Polyadic Decomposition** (NNCPD).
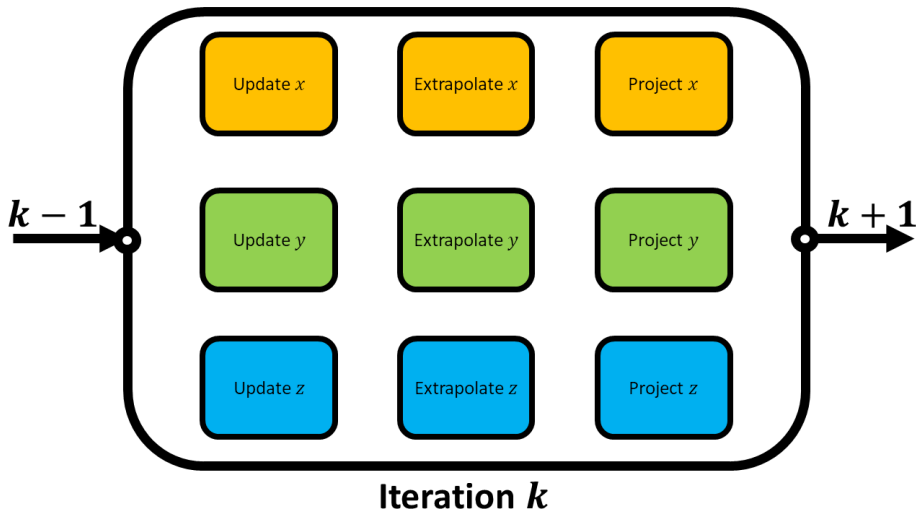
$$\min_{\mathbf{U},\mathbf{V},\mathbf{W}} \Phi(\mathbf{U},\mathbf{V},\mathbf{W}) = \|\mathbf{Y} - \mathbf{U} * \mathbf{V} * \mathbf{W}\| \text{ s.t. } \mathbf{U} \geq 0, \mathbf{V} \geq 0, \mathbf{W} \geq 0$$

$$= \left\|\mathbf{Y} - \sum_{i}^{r} \mathbf{u}_i * \mathbf{v}_i * \mathbf{w}_i\right\|$$



Experiments showed that the approach is very promising and is able to significantly accelerate the NNCPD algorithms.

**Unsolved problem** : NNCPD has even higher variability on the chain structure.

Understanding the relationship between the data structure (rank size, size of each mode) and the chain structure will be crucial.
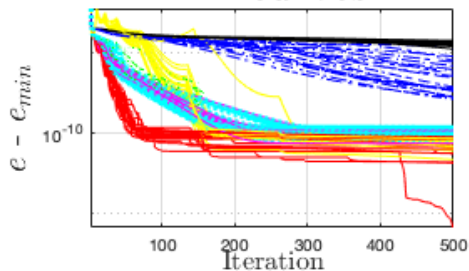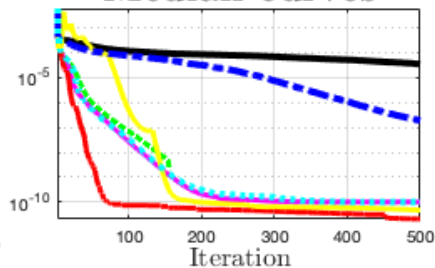
Curves for NNCPD algorithms

What about MU : too slow, not qualified.

TABLE 1 – Median $\mathrm{RE}_{\text{final}}$ in % of $U, V, W$, (* means $\geq 40$)

| Algo | Test 1 | Test 2 | Test 3 |
|------|--------|--------|--------|
| PGD | 15, 1.8, 1.8 | 4.3, *, * | *, *, * |
| PGD-r | 0.2, 1.8, 1.6 | 4.1, *, * | *, *, * |
| HALS | 0.2, 1.6, 1.6 | 2.2, 22, 23 | 4.7, 5.2, 5.2 |
| E-HALS | 0.2, 1.8, 1.8 | 0.04, 0.3, 0.3 | 0.4, 0.8, 0.8 |
| Bro | 0.2, 1.5, 1.5 | 0.2, 2.4, 2.4 | *, *, * |
| APG | 0.5, 3.0, 2.9 | 4.3, *, * | *, *, * |
| APG-r | 0.2, 1.3, 1.2 | 2.7, *, 28 | 10, 11, 11 |

# Outline

## Non-negative Least Square

Problem $(\mathcal{P})$ : given $(\mathbf{A}, \mathbf{b})$, solve

$$\text{(NNLS) } \mathbf{x} = \underset{\mathbf{x} \geq \mathbf{0}}{\operatorname{argmin}} \ \Phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

Let $\mathbf{Q} = \mathbf{A}^\top \mathbf{A}$, $\mathbf{p} = \mathbf{A}^\top \mathbf{b}$, we have an equivalent expression

$$\mathbf{x} = \underset{\mathbf{x} \geq \mathbf{0}}{\operatorname{argmin}} \ \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{p}^\top \mathbf{x} + c$$

$\Phi$ is $\|\mathbf{Q}\|_2$-smooth : the Lipschitz constant of $\nabla\Phi$ is $\|\mathbf{Q}\|_2$

PGD update : $\mathbf{x}^+ = \mathbf{x} - t(\mathbf{Q}\mathbf{x} - \mathbf{p})$ with $t = L^{-1}$

## Multiplicative update

Using the component-wise step size $t_i = \dfrac{x_i}{[\mathbf{Qx}]_i}$, the vector update $\mathbf{x}^+ = \mathbf{x} - t(\mathbf{Qx} - \mathbf{p})$ becomes

$$
\begin{aligned}
x_i^+ &= x_i - t_i([\mathbf{Qx}]_i - p_i) \\
&= x_i - \frac{x_i}{[\mathbf{Qx}]_i}([\mathbf{Qx}]_i - p_i) \\
&= \frac{[\mathbf{Qx}]_i}{[\mathbf{Qx}]_i}x_i - \frac{[\mathbf{Qx}]_i - p_i}{[\mathbf{Qx}]_i}x_i \\
&= \frac{p_i x_i}{[\mathbf{Qx}]_i}
\end{aligned}
$$

In vector form, we have

$$
\mathbf{x}^+ = \mathbf{x} \otimes \frac{\mathbf{p}}{\mathbf{Qx}},
$$

where the multiplication $\otimes$ and division $\frac{[]}{[]}$ are element-wise.

As $\mathbf{p}$, $\mathbf{Q}$ and $\mathbf{x}_0$ are all non-negative, thus the iteration produce a non-negative output.

## Solving NNLS by MU algorithm

Problem :
$$\mathbf{x}_{\mathsf{NNLS}} := \mathrm{argmin}_{\mathbf{x} \geq 0} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

The Multiplicative Update algorithm for NNLS

---
**Algorithm** MU for NNLS
---
**Input:** $\mathbf{A} \in \mathbb{R}_+^{m \times n}$, $b \in \mathbb{R}^m$, an initialization $\mathbf{x} \in \mathbb{R}_+^n$
**Output: x**

1: **for** $k = 1, 2, \ldots$ **do**
2: $\quad \mathbf{x}_{k+1} = \mathbf{x}_k \otimes \dfrac{\mathbf{p}}{\mathbf{Q}\mathbf{x}_k}$
3: **end for**

---

It can be proved that, the objective function $f(\mathbf{x})$ is non-increasing under MU iteration $\mathbf{x}_{k+1} = \mathbf{x}_k \otimes \dfrac{\mathbf{p}}{\mathbf{Q}\mathbf{x}_k}$.

## Solving NNLS by PGD algorithm

Problem :

$$\mathbf{x}_{\text{NNLS}} := \operatorname{argmin}_{\mathbf{x} \geq 0} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

The PGD algorithm for NNLS

---
**Algorithm** PGD for NNLS

---
**Input:** $\mathbf{A} \in \mathbb{R}_+^{m \times n}$, $b \in \mathbb{R}^m$, an initialization $\mathbf{x} \in \mathbb{R}_+^n$
**Output:** $\mathbf{x}$

1: **for** $k = 1, 2, \ldots$ **do**
2: $\quad \mathbf{x}_{k+1} = \left[ \mathbf{x}_k - \frac{1}{L} (\mathbf{Q}\mathbf{x}_k - \mathbf{p}) \right]_+$
3: **end for**

---

It can be proved that, the objective function $f(\mathbf{x})$ is strictly decreasing under PGD iteration when sufficient descent condition holds.

## Solving NNLS by Accelerated-PGD algorithm

Problem :

$$\mathbf{x}_{\mathsf{NNLS}} := \mathrm{argmin}_{\mathbf{x} \geq 0} f(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

The A-PGD algorithm for NNLS

---

**Algorithm** A-PGD for NNLS

---

**Input:** $\mathbf{A} \in \mathbb{R}_+^{m \times n}$, $b \in \mathbb{R}^m$, an initialization $\mathbf{x} \in \mathbb{R}_+^n$
**Output:** $\mathbf{x}$

1: **for** $k = 1, 2, \ldots$ **do**
2:      Compute $\beta_k$
3:      $\mathbf{y}_{k+1} = \left[\mathbf{x}_k - \dfrac{1}{L}(\mathbf{Q}\mathbf{x}_k - \mathbf{p})\right]_+$
4:      $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \beta_k(\mathbf{y}_{k+1} - \mathbf{y}_k)$
5: **end for**

---

# Solving NNLS by Accelerated-PGD algorithm, with restart

Problem :
$$\mathbf{x}_{\mathsf{NNLS}} := \mathrm{argmin}_{\mathbf{x} \geq 0} f(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

The A-PGD algorithm for NNLS

**Algorithm**  A-PGD for NNLS

**Input:** $\mathbf{A} \in \mathbb{R}_+^{m \times n}$, $b \in \mathbb{R}^m$, an initialization $\mathbf{x} \in \mathbb{R}_+^n$
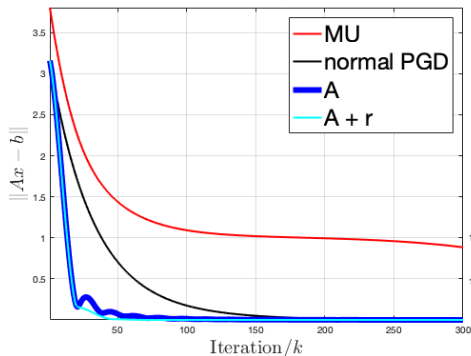**Output: x**

1: **for** $k = 1, 2, \dots$ **do**
2:    Compute $\beta_k$
3:    $\mathbf{y}_{k+1} = \left[\mathbf{x}_k - \dfrac{1}{L}(\mathbf{Q}\mathbf{x}_k - \mathbf{p})\right]_+$
4:    $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \beta_k(\mathbf{y}_{k+1} - \mathbf{y}_k)$
5:    **IF** error increase **do**
6:        $\mathbf{x}_{k+1} = \mathbf{y}_{k+1}$ (take no extrapolation)
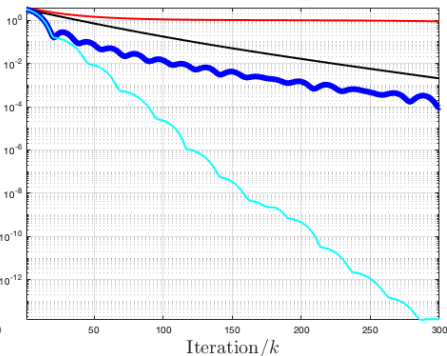7:        reset $\beta$
8:    **ENDIF**
9: **end for**

## Toy example

PGD without any acceleration is already much faster than MU. Not to mention those with acceleration and restart



$(m, n) = 100, 10.$

**What about my scheme?** : With a "good" parameter, the scheme is even faster than Nesterov's type acceleration algorithm. However, all of them are still in linear convergence rate.

You sure want to read it ?
(show the long proof)

# Last page – summary

- What is Non-negative Matrix Factorization, Why NMF
- How to solve NMF *fast* with extrapolation

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

- How to solve NTF *fast* with extrapolation

  **A**.-Cohen-Gillis, "Accelerating Approximate Nonnegative Canonical Polyadic Decomposition using Extrapolation", 2019.

- How to solve NNLS *fast* with extrapolation

  work in progress

- Some open problems

  END OF PRESENTATION.

  Slide, code, preprint in *angms.science*