

Advanced Machine Learning

Lecture 5: Non-Negative Matrix Factorization

Nora Ouzir : nora.ouzir@centralesupelec.fr

Lucca Guardiola : lucca.guardiola@centralesupelec.fr

Oct. - Nov. 2020



CentraleSupélec

Content

1. Reminders on ML
2. Robust regression
3. Hierarchical clustering
4. Classification and supervised learning
5. Non-negative matrix factorization
6. Mixture models fitting
7. Model order selection
8. Dimension reduction and data visualization

Non-negative Matrix Factorization

- ▶ Dimension reduction technique
- ▶ Extensions to deep-NMFs
- ▶ Find a better representation of data to be used for regression, classification, ...

What are the underlying models, how to learn these models, algorithms?

Today's Lecture

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Matrix Factorization: Principles

r = 12 norm

Given a set of data entries $\mathbf{x}_j \in \mathbb{R}^p, 1 \leq j \leq n$ and a **target dimension** $r < \min(p, n)$, we search for r basis elements $\mathbf{w}_k, 1 \leq k \leq r$ and weights $\mathbf{h}_j \in \mathbb{R}^r$ such that

$$\mathbf{x}_j \approx \sum_{k=1}^r \mathbf{w}_k h_j(k)$$

$$\mathbf{X} \approx \mathbf{W}\mathbf{H}$$

Equivalently

Notations

- ▶ $\mathbf{X} \in \mathbb{R}^{p \times n}$ s.t $\mathbf{X}(:, j) = \mathbf{x}_j$ for $1 \leq j \leq n$,
- ▶ $\mathbf{W} \in \mathbb{R}^{p \times r}$ s.t $\mathbf{W}(:, k) = \mathbf{w}_k$ for $1 \leq k \leq r$,
- ▶ $\mathbf{H} \in \mathbb{R}^{r \times n}$ s.t $\mathbf{H}(:, j) = \mathbf{h}_j$ for $1 \leq j \leq n$,

$$\mathbf{X}^{p \times n} = \mathbf{W}^{p \times r} \mathbf{H}^{r \times n}$$

$p \times n$

Matrix Factorization: Key Aspects

Goal: Low rank approximation/ dimension reduction

$$\boxed{X \approx WH}$$

$\begin{pmatrix} p \times r & r \times n \end{pmatrix}$

Key Aspects

1. How to evaluate the quality of the approximation ?
 - ▶ Examples: Frobenius norm, KL-divergence, logistic, Itakura-Saito.
 - Loss function

Matrix Factorization: Key Aspects

Goal: Low rank approximation/ dimension reduction

$$X \approx WH$$

Key Aspects

1. How to evaluate the quality of the approximation ?
 - ▶ Examples: Frobenius norm, KL-divergence, logistic, Itakura-Saito.
 - Loss function
2. Assumptions on the structure of W and H ?
 - ▶ Independence, sparsity, normalization, ...
 - Non-negativity

NMF: find (W, H) s.t $X \approx WH$, $W \geq 0, H \geq 0$

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

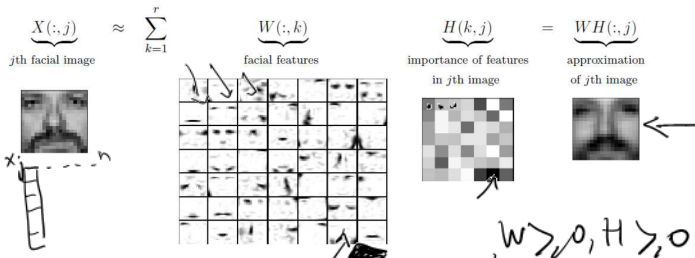
2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Example: Facial Feature Extraction

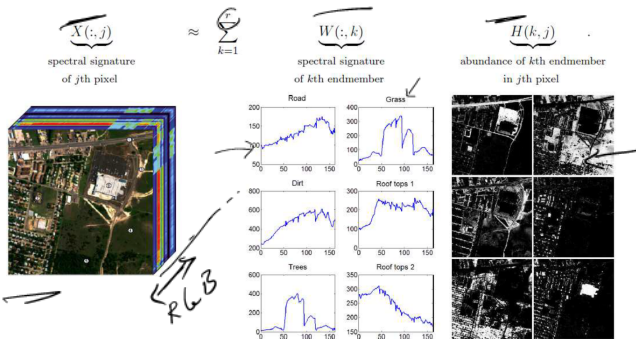


Decomposition of the CBCL face database [Lee and Seung, 1999]

- ▶ Some of the features look like parts of facial features/organs: noses, eyes, etc.
- ▶ Decomposition of a face as having a certain weight of a some nose type, a certain amount of some eye type, etc.

→ Interpretable representation

Example: Spectral Unmixing




Decomposition of the Urban hyperspectral image [Ma et al., 2014]

- NMF is able to compute the spectral signatures of the endmembers and simultaneously the abundance of each endmember in each pixel.

→ Here non-negativity is directly relevant for the application

Example: Topic Modelling in Text Mining

Goal: Decompose a term-document matrix, where each column represents a document, and each element in the document represents the weight of a certain word (e.g., term frequency - inverse document frequency).



$$\rightarrow \underbrace{X(:, j)}_{\text{jth document}} \approx \sum_{k=1}^r \underbrace{W(:, k)}_{\text{kth topic}} \underbrace{H(k, j)}_{\text{importance of kth topic in jth document}}$$

Topic decomposition model [Blei, 2012]

- ▶ The ordering of the words in the documents is not taken into account (=bag-of-words).
- ▶ The NMF decomposition of the term-document matrix yields components that could be considered as "topics"
- ▶ Decomposes each document into a weighted sum of topics.

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

White board

- $X \approx WH$
- $(X, WH) \rightarrow \underline{D}(X, WH) = \sum_{i=1}^n \sum_{j=1}^p d(x_{ij}, [WH]_{ij})$
 - \rightarrow Separable
 - $\rightarrow \begin{cases} d(x, x) = 0 \\ d(x, y) \geq 0 \end{cases}$
- Example $d(x, y) = (x - y)^2$
QUADRATIC

$\begin{aligned} W^T W &= Id \\ H^T H &= \text{diag}(\lambda) \end{aligned}$ <p>SVD</p>	$\begin{aligned} H^T H &= Id \\ H_{kj} &\in \{0, 1\} \end{aligned}$ <p>k-MEANS</p>	$\begin{aligned} W &\geq 0 \\ H &\geq 0 \end{aligned}$ <p>LS-NMF</p>
---	--	--

• Gaussian \rightarrow nonrealistic
 \rightarrow robust

White board

✗ QUADRATIC.

✗ KL-divergence

Itakura-Saito $\rightarrow d(x, y) = \frac{x}{y} - \log \frac{x}{y} - 1$

Music Analysis $d(x, y) = d(y, x)$

\rightarrow Need more advanced solvers

$$d(x, y) = x \log \frac{x}{y} + x - y$$

$1/w \geq 0, H \geq 0$ KL NMF.

White board

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

NMF Algorithms

Recall: we want to find $(\underline{W}, \underline{H})$ s.t. $X \approx WH$, $\underline{W} \geq 0$, $\underline{H} \geq 0$

- ▶ Loss function + constraints (+ regularization)
- Non-convex optimization problem wrt \underline{W} and \underline{H}

Main types of algorithms

- ▶ Multiplicative methods
- ▶ Alternating LS
- ▶ Gradient descent

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Multiplicative Algorithms for NMF

Solve

$$\min_{W, H} \mathcal{D}(X, WH) \text{ s.t. } W \geq 0, H \geq 0$$

Challenges

NMF is **NP-hard** and **ill-posed** (solution is not unique):

- ▶ Most algorithms are only guaranteed to converge to a stationary point
- ▶ Sensitive to initialization
- ▶ In practice: priors on W and H /regularization

[Lee and Seung, 1999]

- Popular class of methods relying on multiplicative updates
- Key assumption $X \geq 0$

Multiplicative Algorithms: [Lee and Seung]

Assuming $X \geq 0$, the resulting updates are multiplicative and take the following forms

- Frobenius norm $\|X - WH\|_F^2$

elementwise

$$W \leftarrow W \circ \frac{XH^T}{WHH^T}$$

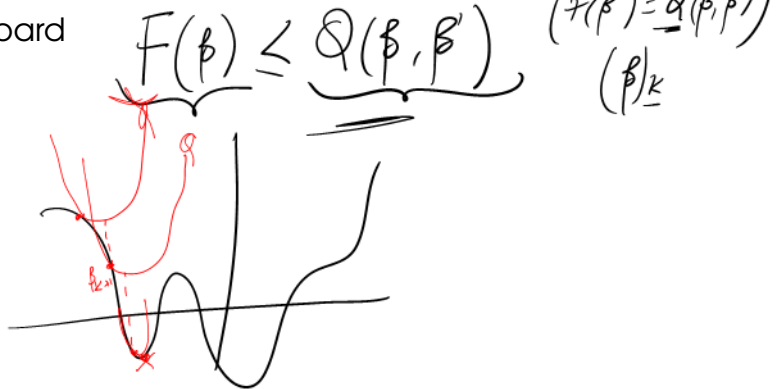
$$H \leftarrow H \circ \frac{W^T X}{W^T W H}$$

- KL-divergence: $\mathcal{KL}(X, WH)$

$$W_{ik} \leftarrow W_{ik} \frac{\sum_{l=1}^n (H_{kl} X_{il} / [WH]_{il})}{\sum_{l=1}^n H_{kl}}$$

$$H_{kj} \leftarrow H_{kj} \frac{\sum_{i=1}^p (W_{ik} X_{ij} / [WH]_{ij})}{\sum_{i=1}^p W_{ik}}$$

White Board



USE & STUNG : using Quadratic Loss.

White Board

FROBENIUS: $\|WH - X\|_F^2 = \sum_{j=1}^L \|Wh_j - x_j\|^2$

- 1 Build Q
- 2 Separable w.r.t lines (i)

$\beta' = \bar{H} \bar{H}^T$ $\|Wh_j - x_j\|^2 \leq \sum_{i=1}^P \frac{1}{[w \bar{h}_j]_i} \sum_{k=1}^r w_{ik} \bar{H}_{kj} (x_{ij} - \frac{H_{kj}}{\bar{H}_{kj}} [w \bar{h}_j]_i)^2$

JENSEN INEQUALITY

$f(\sum \lambda_i x_i) \leq \sum \lambda_i f(x_i)$

can show that $\bar{H} = H$

$F \leq \sum_{j=1}^L \sum_{i=1}^P \sum_{k=1}^r \frac{w_{ik} \bar{H}_{kj}}{[w \bar{h}_j]_i} (x_{ij} - \frac{H_{kj}}{\bar{H}_{kj}} [w \bar{h}_j]_i)^2$

$Q(H, \bar{H})$

White Board

$$\frac{\partial Q}{\partial H_{kj}} = 2 \sum_i \frac{w_{ik} \overline{H_{kj}}}{[\overline{w h_j}]_i} \times \sum \frac{[\overline{w h_j}]_i}{\overline{H_{kj}}} \times \left(x_{ij} - \frac{H_{kj}}{\overline{H_{kj}}} [\overline{w h_j}]_i \right)$$

$$\frac{\partial Q}{\partial H_{kj}} = 0 \Rightarrow \sum_{i=1}^P w_{ik} \left(x_{ij} - \frac{H_{kj}}{\overline{H_{kj}}} [\overline{w h_j}]_i \right) = 0$$

$$\Rightarrow \sum_{i=1}^P \underbrace{w_{ik} x_{ij}}_{[W^T X]_{kj}} = \frac{H_{kj}}{\overline{H_{kj}}} \sum_{i=1}^P \underbrace{w_{ik} [\overline{w h_j}]_i}_{[W^T W \overline{H}]_{kj}}$$

$$x_{ij} > 0$$

$$H_{kj} = \overline{H_{kj}} \times \frac{[W^T X]_{kj}}{[W^T W \overline{H}]_{kj}}$$

$$\rightarrow Q(H, \overline{H})$$

minimize

$$\rightarrow \frac{\partial Q}{\partial H_{kj}} = 0 \rightarrow MU$$

Multiplicative Algorithms: Surrogate functions

The multiplicative schemes rely on **separable surrogate functions** that **majorize** the loss w.r.t \mathbf{W} and \mathbf{H} . For every $(\mathbf{X}, \mathbf{W}, \mathbf{H}, \bar{\mathbf{H}}) \geq 0$, and $1 \leq j \leq n$:

► **Frobenius norm:**

$$\|\mathbf{W}\mathbf{h}_j - \mathbf{x}_j\|_2^2 \leq \sum_{i=1}^p \frac{1}{[\mathbf{W}\bar{\mathbf{h}}_j]_i} \sum_{k=1}^r w_{ik} \bar{H}_{kj} \left(x_{ij} - \frac{H_{kj}}{\bar{H}_{kj}} [\mathbf{W}\bar{\mathbf{h}}_j]_i \right)^2$$

► **KL-divergence:**

$$\begin{aligned} \mathcal{KL}(\mathbf{x}_j, \mathbf{W}\mathbf{h}_j) &\leq \sum_{i=1}^p \left[X_{ij} \log X_{ij} - X_{ij} + [\mathbf{W}\bar{\mathbf{h}}_j]_i \right. \\ &\quad \left. - \frac{X_{ij}}{[\mathbf{W}\bar{\mathbf{h}}_j]_i} \sum_{k=1}^r w_{ik} \bar{H}_{kj} \log \left(\frac{H_{kj}}{\bar{H}_{kj}} [\mathbf{W}\bar{\mathbf{h}}_j]_i \right) \right] \end{aligned}$$

Multiplicative Algorithms: Implementation Key points

MU algorithms

- ▶ Build a quadratic majorant function
- ▶ Jensen inequality → separability
- ▶ Guarantee positive updates by construction

Implementation: initialization

- ▶ Apply factorization without non-negativity constraint (e.g., SVD)
→ truncate/threshold
- ▶ Avoid too simple initializations (e.g., identity, random)
- ▶ Clustering methods, r columns of X ...

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Weighted NMF

Application example: Matrix completion to predict unobserved data (e.g., user-rating) → Binary weights indicate the position of the available entries in X .

- Weighted Frobenius norm $\|\Sigma \circ (X - WH)\|_F^2$



$$W \leftarrow W \circ \frac{(\Sigma \circ X)H^T}{(\Sigma \circ WH)H^T}$$

$$H \leftarrow H \circ \frac{W^T(\Sigma \circ X)}{W^T(\Sigma \circ WH)}$$

- Weighted KL-divergence: $\mathcal{KL}(X, \text{Diag}(p)WH\text{Diag}(q))$

$p > 0$
 $q > 0$

$$W_{ik} \leftarrow W_{ik} \frac{\sum_{l=1}^n (H_{kl} X_{il} / \textcircled{p}_i) [WH]_{il}}{\sum_{l=1}^n \underline{q}_l H_{kl}}$$

$$H_{kj} \leftarrow H_{kj} \frac{\sum_{i=1}^p (W_{ik} X_{ij} / \textcircled{q}_j) [WH]_{ij}}{\sum_{i=1}^p \textcircled{p}_i W_{ik}}$$

White Board

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Regularized NMFs

Regularized Frobenius norm: (sparsity, stability, scale ambiguity)

$$\frac{1}{2} \|X - WH\|_F^2 + \frac{\mu}{2} \|H\|_F^2 + \lambda \|H\|_1 + \frac{\nu}{2} \|W\|_F^2 - \dots$$

Resulting MU

$$W \leftarrow W \circ \frac{XH^T}{W(HH^T + \nu I_r)}$$

$$H \leftarrow H \circ \frac{W^T X - \lambda \mathbf{1}_{r \times n}}{(W^T W + \mu I_r)H}$$

\sum_h $L_1 \rightarrow \text{sparsity}$

- The regularization terms are already separable!
- The ambiguity due to the rescaling and rotation of (W, H) is frozen by the penalty terms.

White Board

Multiplicative NMF Algorithms: Summary

To sum up

- ▶ Dimension reduction for non-negative data
- ▶ Easily interpretable representation
- ▶ Weighted and regularized formulations

In practice MU is...

- + Simple to implement
- Slow convergence
- Sensitive to initialization (stuck with small values)
- Choice of r ? Heuristics exist, domain knowledge,...

Today's course

1. Introduction

1. Main principles of NMF
2. Key applications

2. Common Losses for NMF

3. Multiplicative Algorithms

1. Quadratic and KL distances
2. Weighted NMFs
3. Regularized NMFs

4. Other NMF Algorithms

Other NMF Algorithms: Alternating LS

The subproblem in one variable is convex!

Alternating LS: Unconstrained solution w.r.t. W or H followed by a projection onto non-negative orthant.

- ▶ Easy to implement but oscillations can arise (no convergence guarantees) → initialization purposes.

Other NMF Algorithms: Alternating LS

The subproblem in one variable is convex!

Alternating LS: Unconstrained solution w.r.t. W or H followed by a projection onto non-negative orthant.

- ▶ Easy to implement but oscillations can arise (no convergence guarantees) → initialization purposes.

Alternating Non-negative LS: Solve the constrained problem (alternate w.r.t. W, H) exactly using an inner solver (e.g., projected gradient, Quasi-Newton, active set).

- ▶ Expensive → useful as a refinement step for a cheaper MU.

Other NMF Algorithms: Alternating LS

The subproblem in one variable is convex!

Alternating LS: Unconstrained solution w.r.t. W or H followed by a projection onto non-negative orthant.

- ▶ Easy to implement but oscillations can arise (no convergence guarantees) → initialization purposes.

Alternating Non-negative LS: Solve the constrained problem (alternate w.r.t. W, H) exactly using an inner solver (e.g., projected gradient, Quasi-Newton, active set).

- ▶ Expensive → useful as a refinement step for a cheaper MU.

Hierarchical Alternative LS: Exact coordinate descent method, updating one column of W (resp. one line of H) at a time.

- ▶ Simple to implement, performance is similar to MU.

Other NMF Algorithms: Alternate Scheme

Algorithm CD Two-Block Coordinate Descent – Framework of Most NMF Algorithms

Input: Input nonnegative matrix $X \in \mathbb{R}_+^{p \times n}$ and factorization rank r .

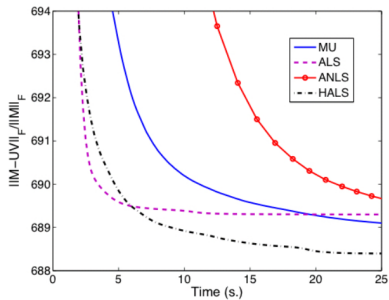
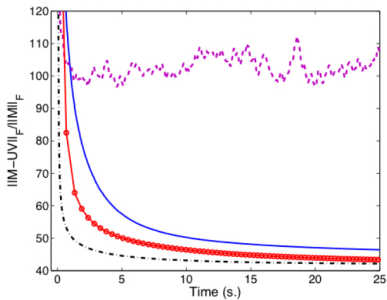
Output: $(W, H) \geq 0$: A rank- r NMF of $X \approx WH$.

- 1: Generate some initial matrices $W^{(0)} \geq 0$ and $H^{(0)} \geq 0$; see Section 3.1.8
- 2: **for** $t = 1, 2, \dots^\dagger$ **do**
- 3: $W^{(t)} = \text{update}(X, H^{(t-1)}, W^{(t-1)})$.
- 4: $H^{(t)T} = \text{update}(X^T, W^{(t)T}, H^{(t-1)T})$.
- 5: **end for**

† See Section 3.1.7 for stopping criteria.

Algorithm: Coordinate descent framework

Other NMF Algorithms: Comparison



Performance comparison to MU using same initialization:
Average error vs execution times