# Chapter 2
# **Non-Negative Matrix Factorization**

Part 1: Introduction & computation

max planck institut
informatik

# **Motivating NMF**

Pauli Miettinen

# Reminder

$$\boldsymbol{A}$$

| 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 |

**=**

$$\boldsymbol{W}_1\boldsymbol{\Sigma}_{1,1}\boldsymbol{V}_1^T$$

| 0.6 | 1.3 | 0.6 | 1.3 | 0.6 |
|-----|-----|-----|-----|-----|
| 0.3 | 0.8 | 0.3 | 0.8 | 0.3 |
| 0.3 | 0.8 | 0.3 | 0.8 | 0.3 |

| 0.3 |
|-----|
| 0.5 |

**+**

$$\boldsymbol{W}_2\boldsymbol{\Sigma}_{2,2}\boldsymbol{V}_2^T$$

| 0.4 | −0.3 | 0.4 | −0.3 | 0.4 |
|-----|------|-----|------|-----|
| −0.3 | 0.2 | −0.3 | 0.2 | −0.3 |
| −0.3 | 0.2 | −0.3 | 0.2 | −0.3 |

The components of the SVD are not very interpretable

# Non-negative factors

$$A$$

| 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 |

=

$W$   $W_1 H_1$       $H$         $W_2 H_2$

| 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |

+

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 |

Forcing the factors to be non-negative can, and often will, improve the interpretability of the factorization

# The definition

# Definition of NMF

Given a non-negative matrix $A \in \mathbb{R}_+^{n \times m}$ and integer $k$, find non-negative matrices $W \in \mathbb{R}_+^{n \times k}$ and $H \in \mathbb{R}_+^{k \times m}$ such that
$$\frac{1}{2}\|A - WH\|_F^2$$
is minimized.

# Non-negative rank

- The **non-negative rank** of matrix $A$, $\text{rank}_+(A)$, is the size of the smallest exact non-negative factorization $A = WH$

  - $\text{rank}(A) \leq \text{rank}_+(A) \leq \min\{n, m\}$

# Some comments

- NMF is **not** unique

  - If $X$ is nonnegative and with nonnegative inverse, then $WXX^{-1}H$ is equivalent valid decomposition

- Computing NMF (and non-negative rank) is NP-hard

  - This was open until 2008

# Example of non-uniqueness

# NMF has no order

- The factors in NMF have no inherent order

  - The first component is no more important than the second is no more important…

- NMF is not **hierarchical**

  - The factors of rank-($k$+1) decomposition can be completely different to those of rank-$k$ decomposition

# Example

# Interpreting NMF

# Parts-of-whole

- NMF works over **anti-negative semiring**

  - There is no subtraction

- Each rank-1 component $\boldsymbol{w}_i\boldsymbol{h}_i$ explains a part of the whole

  - This can yield to sparse factors
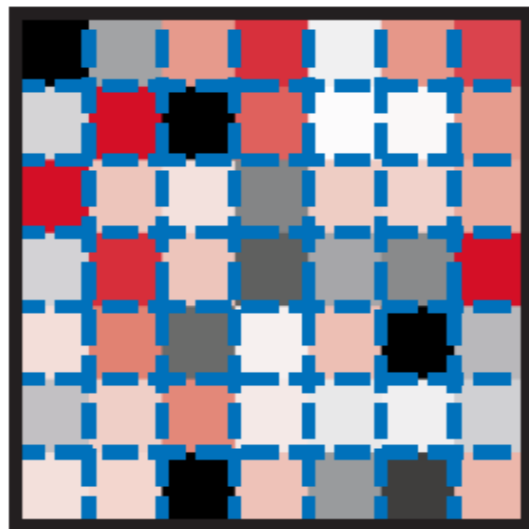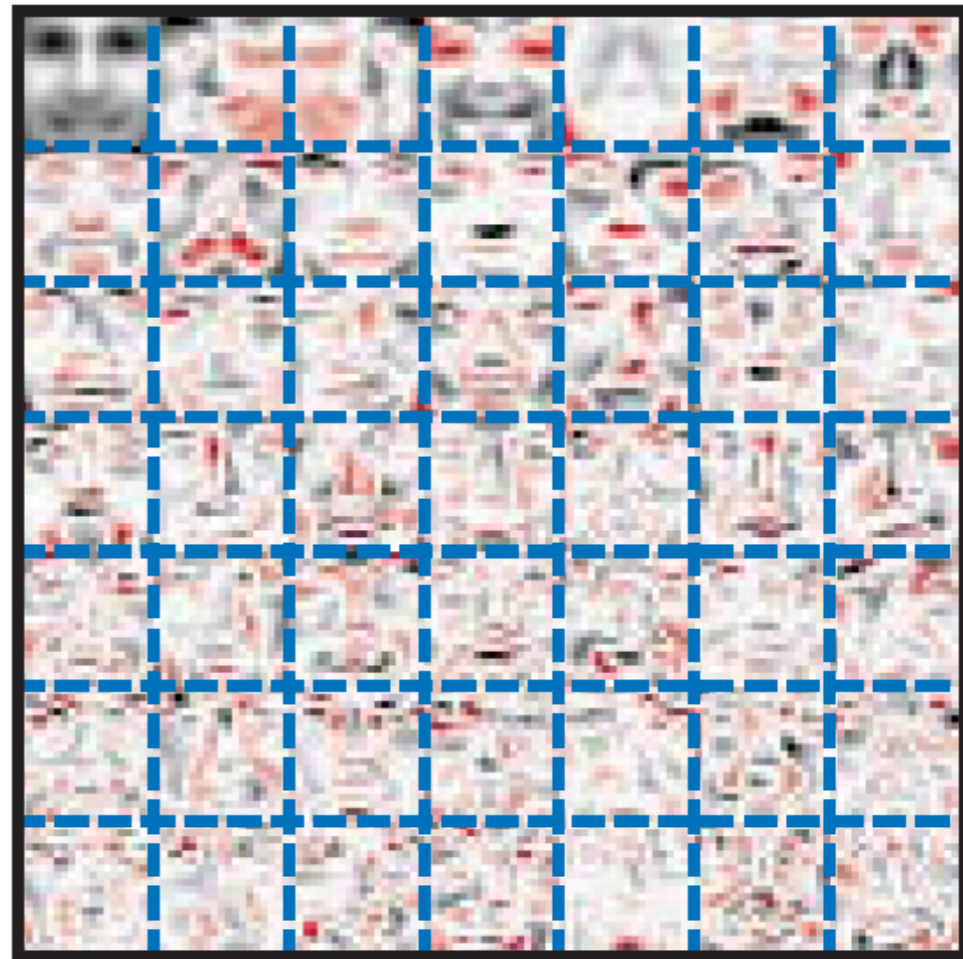
# NMF example: faces



Row of original

Row of reconstruction

PCA/SVD

=   ×

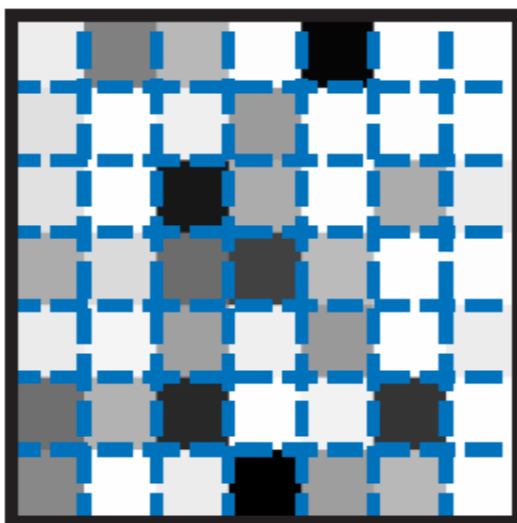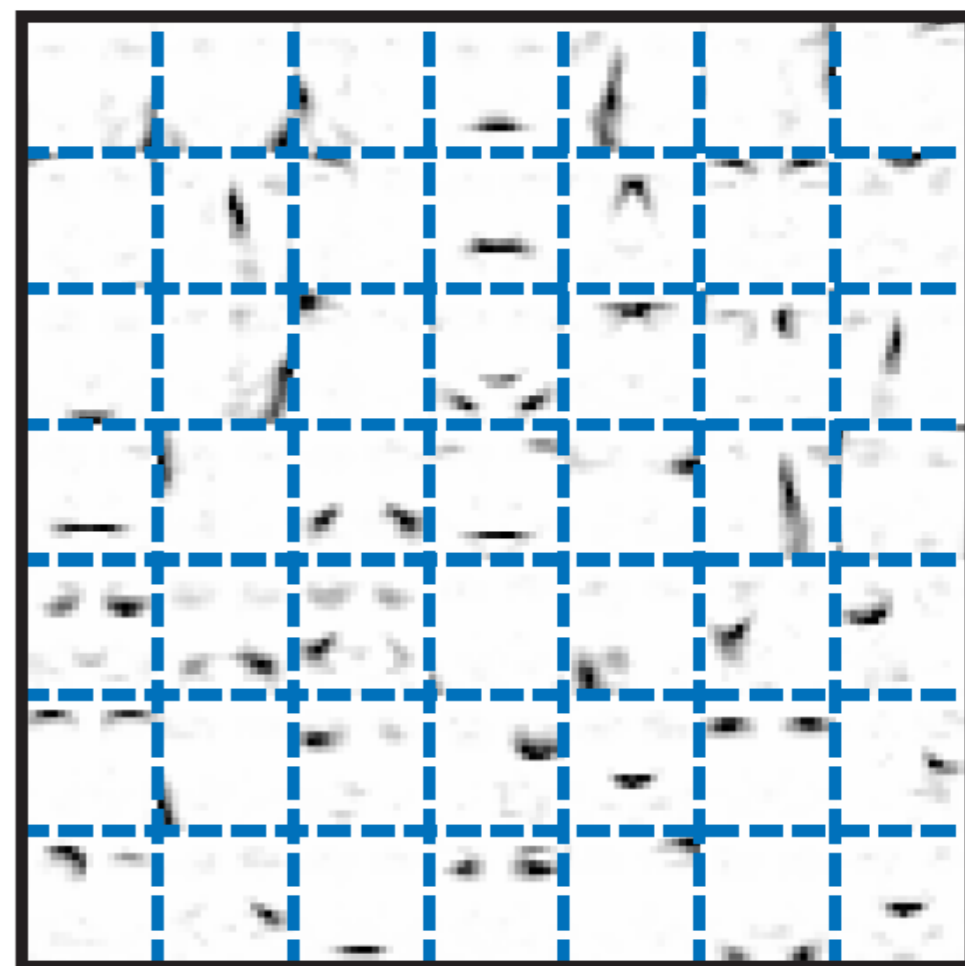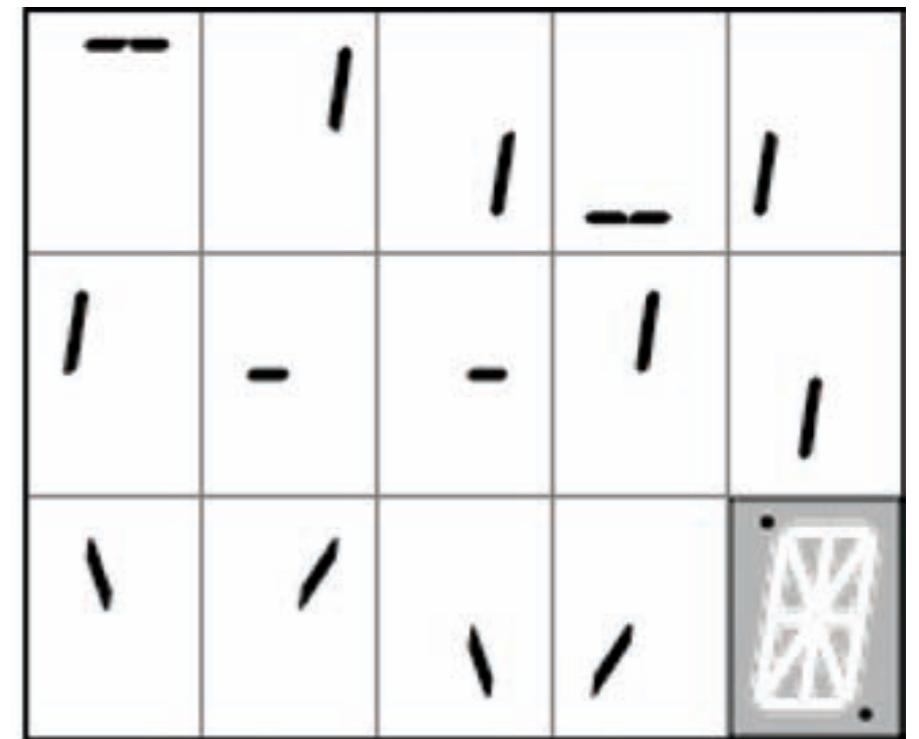# NMF example: faces



Row of original

NMF

# NMF example: digits

NMF factors correspond to patterns and background



*A*

*H*

# Some NMF applications

- Text mining (more later)

- Bioinformatics

- Microarray analysis

- Mineral exploration

- Neuroscience

- Image understanding

- Air pollution research

- Weather forecasting

- …



(a) Original 10 sources

(b) Typical 1000 mixtures

(c) Ten estimated components by using Fast-HALS

(d) PSNR using Beta HALS for various values of $\beta$

**Figure 4.8** Illustration for (a) benchmark used in large-scale experiments with 10 nonnegative sources; (b) Typical 1000 mixtures; (c) Ten estimated components by using FAST HALS NMF from the observations matrix **Y** of dimension $1000 \times 1000$. (d) Performance expressed via the PSNR using the Beta HALS NMF algorithm for $\beta = 0.5, 1, 1.5, 2$ and 3.

# Computing NMF

# General idea

- NMF is not convex, but it is **biconvex**

  - If $W$ is fixed, $\frac{1}{2}\|A - WH\|_F^2$ is convex

- Start from random $W$ and **repeat**

  - Fix $W$ and update $H$

  - Fix $H$ and update $W$

- **until** the error doesn't decrease anymore

# Notes on the general idea

- How to create a good random starting point?

  - Is the algorithm robust to initial solutions?

- How to update **W** and **H**?

- When (and how quickly) has the process converged?

  - Fixed number of iterations? Minimum change in error?

# Alternating least squares

- Without the non-negativity constraint, this is the standard least-squares:

  - $w_i \leftarrow \text{argmin}_w \ ||wH - a_i||_F$

  - we can update $W \leftarrow AH^+$ and $H \leftarrow W^+A$

  - $X^+$ is the pseudo-inverse of $X$ which is LS-optimal

- The method is called **alternating least-squares** (ALS)

- This can introduce negative values

# Enforcing non-negativity in ALS

- Least-squares optimal update of **W** (or **H**) with non-negativity constraints is convex optimization problem

  - In theory in P, in practice slow, but subject to much research

- Simple approach: truncate all negative values to 0

  - Update $\boldsymbol{W} \leftarrow [\boldsymbol{A}\boldsymbol{H}^+]_+$

# The NMF-ALS algorithm

1. $W \leftarrow \text{random}(n, k)$

2. **repeat**

   2.1. $H \leftarrow [W^+ A]_+$

   2.2. $W \leftarrow [A H^+]_+$

3. **until** convergence

# When has there been enough convergence?

- When the error doesn't change too much

  - $||\boldsymbol{A} - \boldsymbol{W}^{(k)}\boldsymbol{H}^{(k)}||_F - ||\boldsymbol{A} - \boldsymbol{W}^{(k+1)}\boldsymbol{H}^{(k+1)}||_F \leq \epsilon$

- After some number of maximum iterations has been achieved
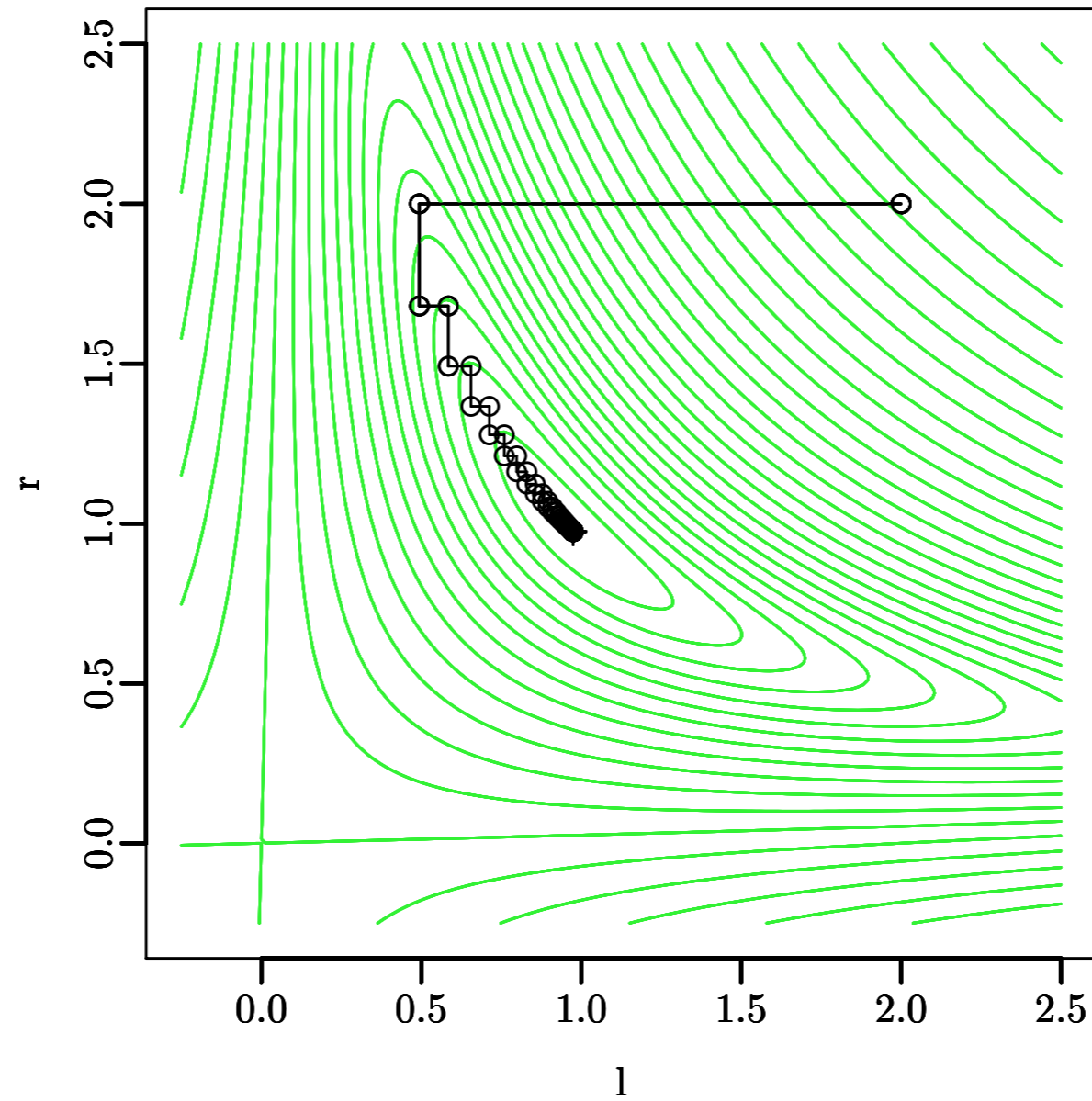
- Usually, whichever of these two happens first

# Gradient descent

- We can compute the gradient of the error function (with one factor matrix fixed)

  - $f(\boldsymbol{H}) = \frac{1}{2} \|\boldsymbol{A} - \boldsymbol{W}\boldsymbol{H}\|_F^2 = \frac{1}{2} \sum_i \|\boldsymbol{a}_i - \boldsymbol{W}\boldsymbol{h}_i\|_F^2$

  - $\nabla_{\boldsymbol{H}_{ij}} f(\boldsymbol{H}) = (\boldsymbol{W}^T\boldsymbol{A})_{ij} - (\boldsymbol{W}^T\boldsymbol{W}\boldsymbol{H})_{ij}$

- We can move slightly towards the negative gradient

  - How much is the step size and deciding it is a big problem

# The NMF gradient descent algorithm

1. $W \leftarrow$ random($n, k$)

2. $H \leftarrow$ random($k, m$)

3. **repeat**

    3.1. $H \leftarrow H - \varepsilon_H \frac{\partial f}{\partial H}$

    3.2. $W \leftarrow W - \varepsilon_W \frac{\partial f}{\partial W}$

4. **until** convergence

# Example

# Notes on gradient descent

- Choosing the correct step size is crucial

  - Usually the shorter step sizes the closer the solution we are

- Can converge to local minimum

  - Wrong step size, and converges very close to the initial solution

# The NMF multiplicative updates algorithm

1. $W \leftarrow$ random($n$, $k$)

2. $H \leftarrow$ random($k$, $m$)

3. **repeat**

   3.1. $h_{ij} \leftarrow h_{ij} \dfrac{(\boldsymbol{W}^T \boldsymbol{A})_{ij}}{(\boldsymbol{W}^T \boldsymbol{W} \boldsymbol{H})_{ij} + \varepsilon}$

   3.2. $w_{ij} \leftarrow w_{ij} \dfrac{(\boldsymbol{A} \boldsymbol{H}^T)_{ij}}{(\boldsymbol{W} \boldsymbol{H} \boldsymbol{H}^T)_{ij} + \varepsilon}$

4. **until** convergence

# Notes on multiplicative updates

- Proposed by Lee & Seung (Nature, 1999)

- Equivalent to gradient descent with dynamic step size

- Zeros in initial solutions will never turn into non-zeros; non-zeros will never turn into zeros

  - Problems if the correct solution contains zeros