

# Chapter 3

# Non-Negative Matrix Factorization

Part 1: Introduction & computation



# Motivating NMF

# Reminder

$$\mathbf{A} = \mathbf{U}_1 \mathbf{\Sigma}_{1,1} \mathbf{V}_1^T + \mathbf{U}_2 \mathbf{\Sigma}_{2,2} \mathbf{V}_2^T$$

1	1	1	1	1
0	1	0	1	0
0	1	0	1	0

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	0.4	-0.3	0.4	-0.3	0.4
0.5	-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

$$\begin{array}{c} \mathbf{A} \\ \begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 1 & 1 & 1 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline \end{array} \end{array} = \begin{array}{c} \mathbf{W} \mathbf{W}_1 \mathbf{H}_1 \quad \mathbf{H} \quad \mathbf{W}_2 \mathbf{H}_2 \\ \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ \hline \end{array} \end{array}$$

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  $\mathbf{A} \in \mathbb{R}_+^{n \times m}$  and an integer  $k$ , find non-negative matrices  $\mathbf{W} \in \mathbb{R}_+^{n \times k}$  and  $\mathbf{H} \in \mathbb{R}_+^{k \times m}$  such that

$$\frac{1}{2} \|\mathbf{A} - \mathbf{WH}\|_F^2$$

is minimized.

# Non-negative rank

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

# 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

$$\begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 1 & 1 & 1 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 1 & 1 & 1 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline \end{array}$$

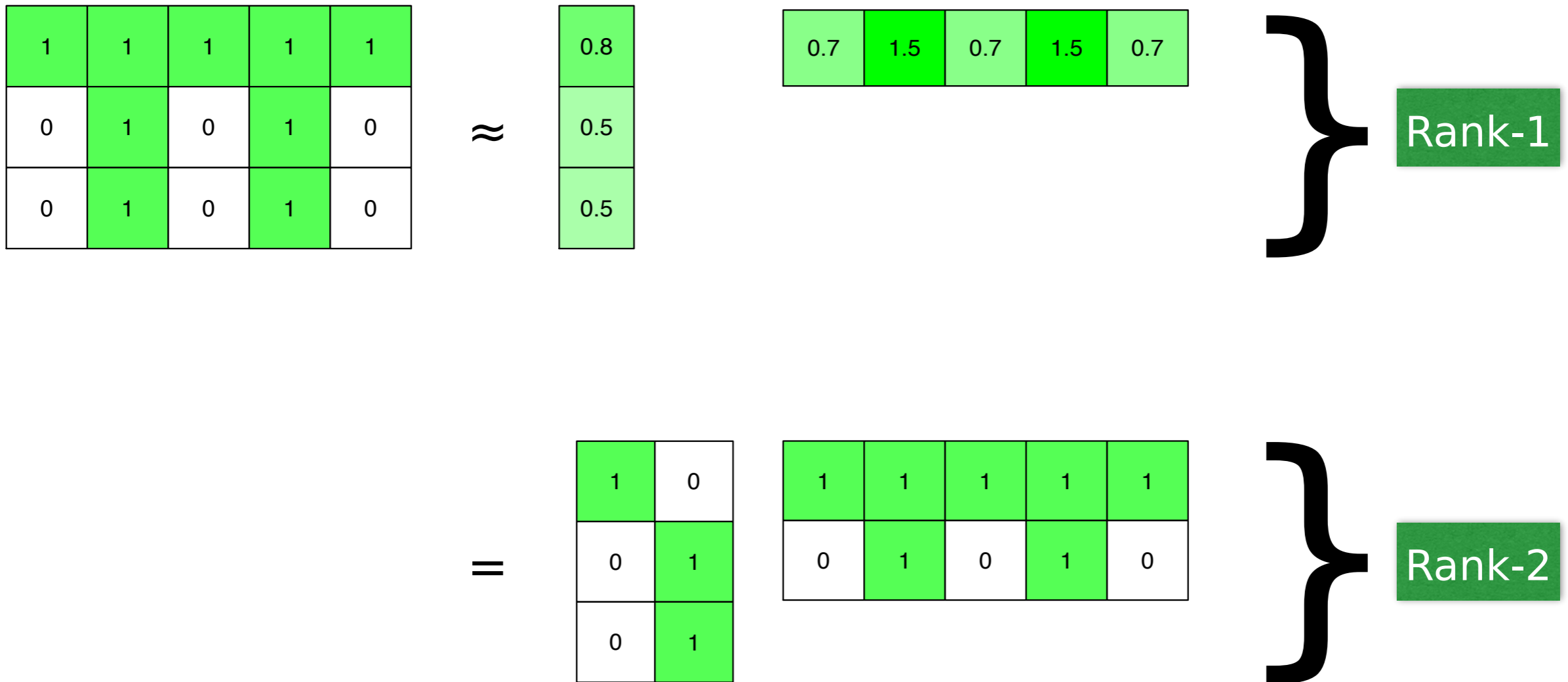
$$= \begin{array}{|c|c|c|c|c|} \hline 1 & 0.5 & 1 & 0.5 & 1 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|c|} \hline 0 & 0.5 & 0 & 0.5 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|c|c|} \hline 1 & 0 & 1 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 1 & 0 \\ \hline \end{array}$$

# 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  $\mathbf{w}_i \mathbf{h}_i$  explains a part of the whole
  - This can yield to sparse factors

# NMF example: faces



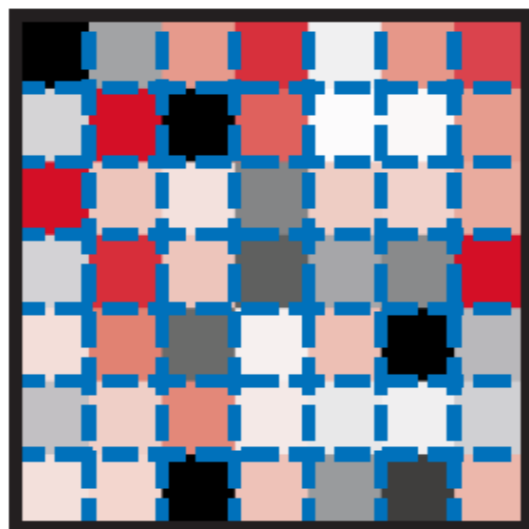
Row of original

PCA/SVD

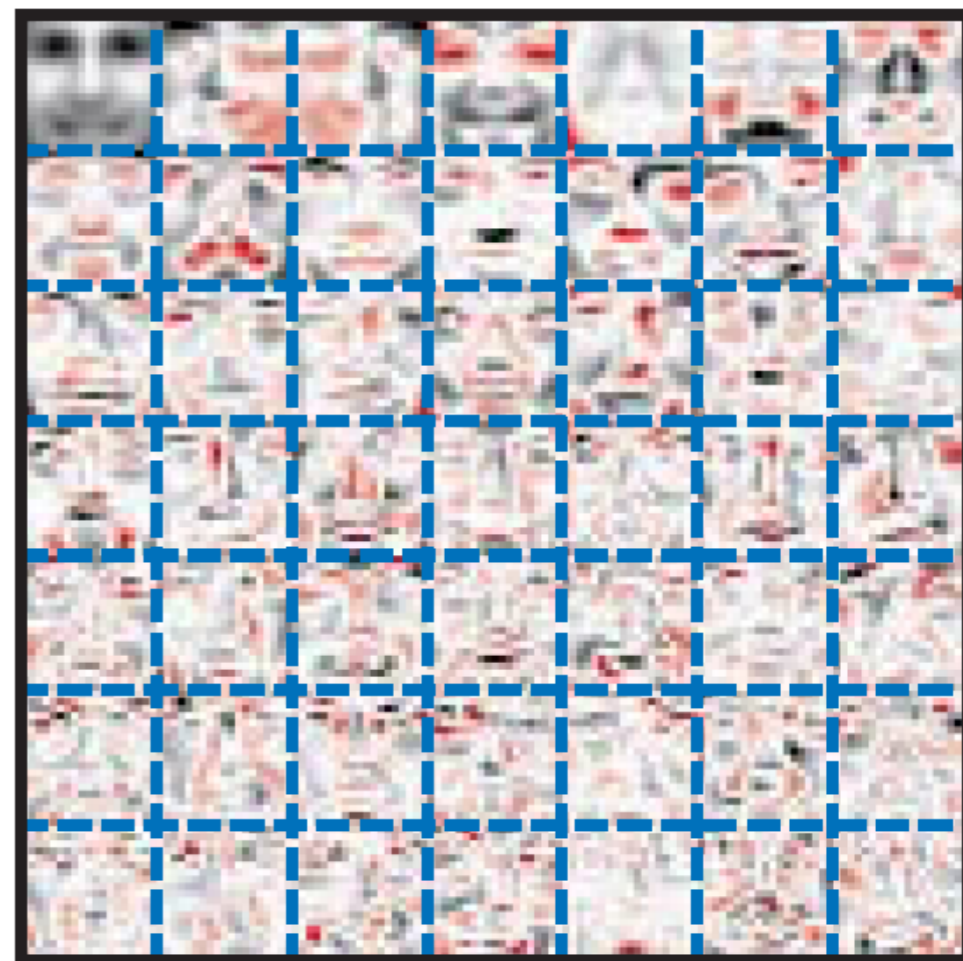


Row of reconstruction

=



×



# NMF example: faces

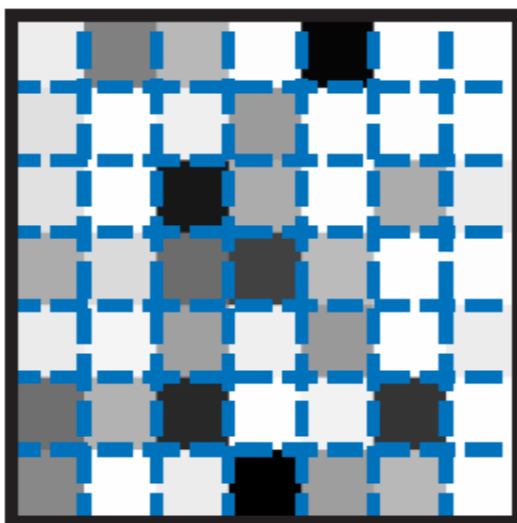


Row of original

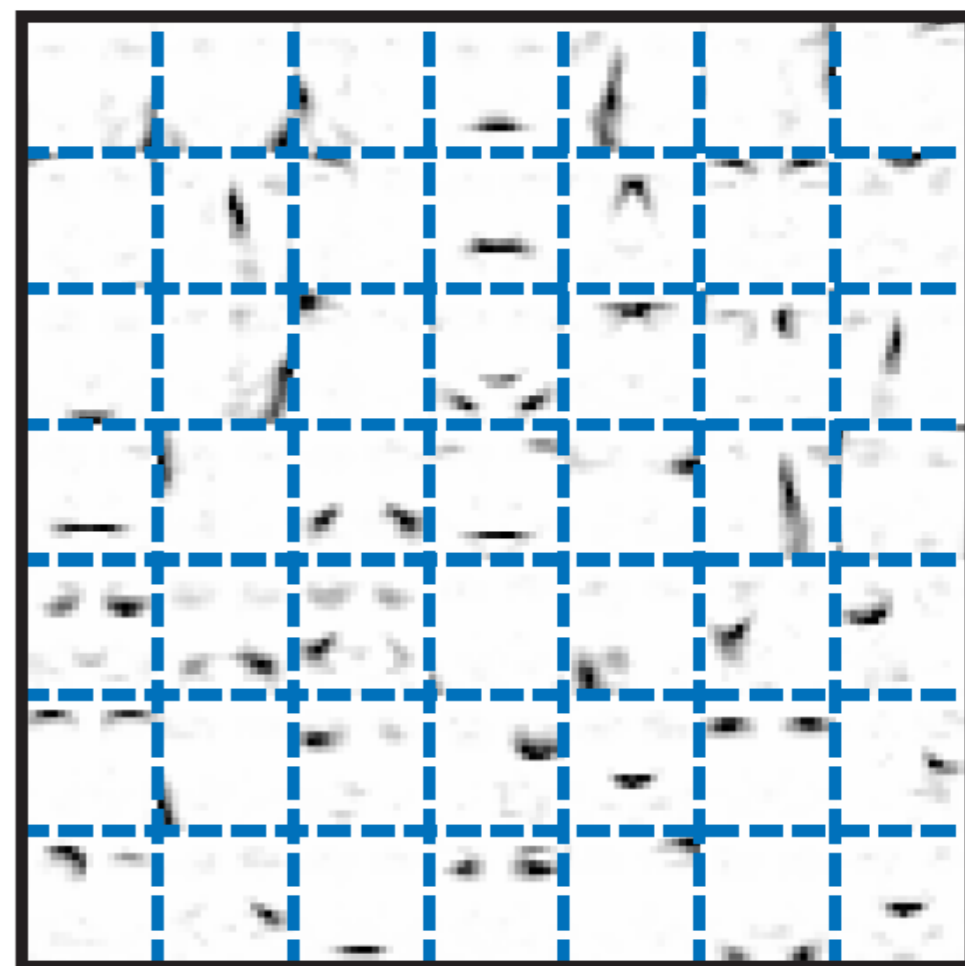
NMF



=



×

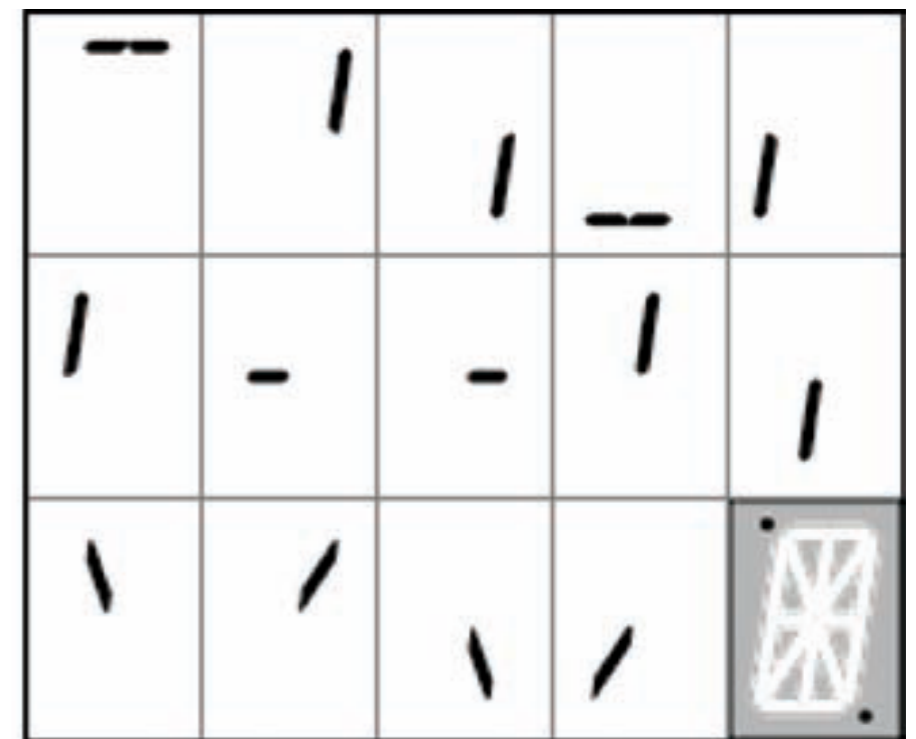


# NMF example: digits

NMF factors correspond to patterns and background



**A**

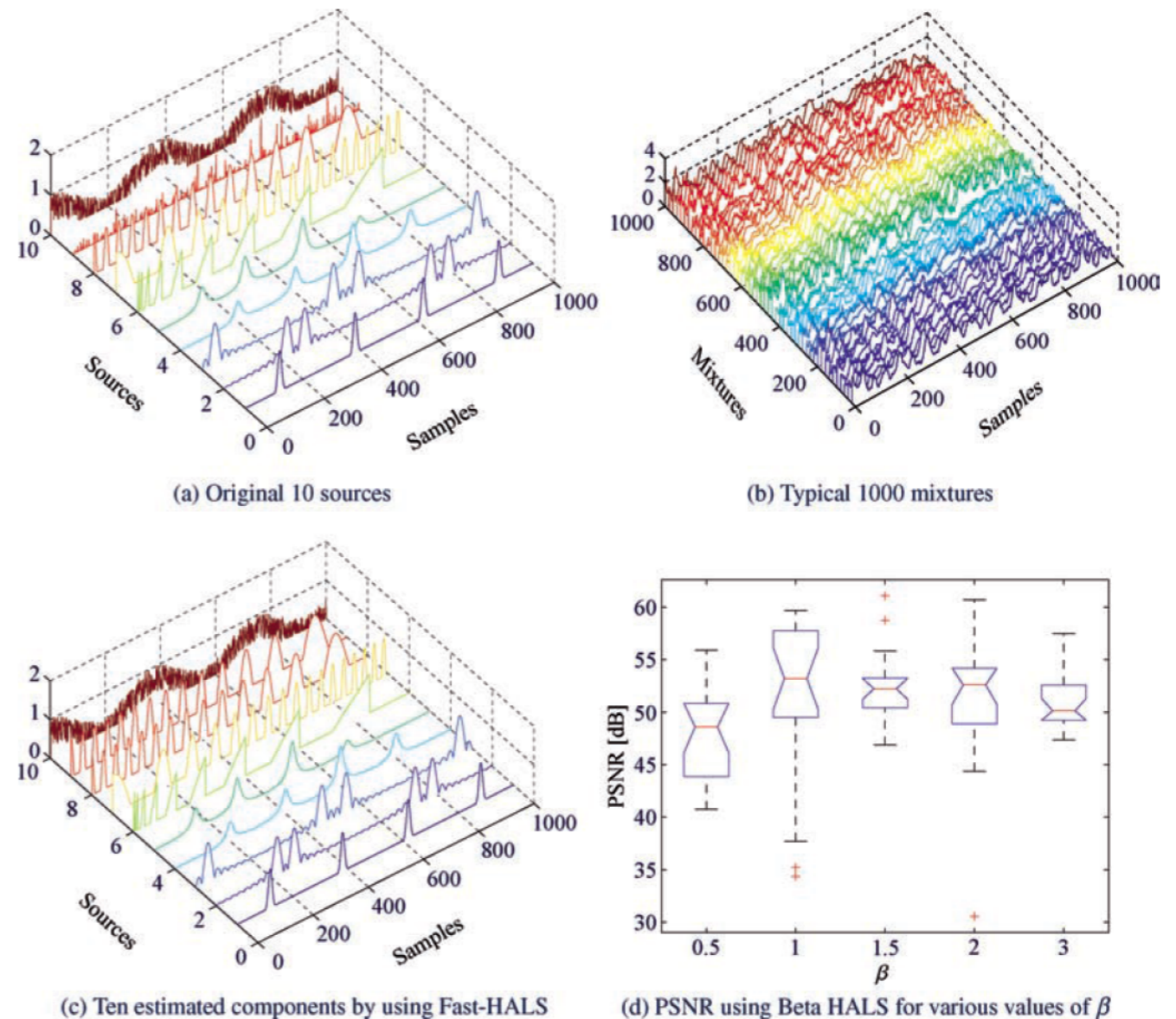


**H**



# Some NMF applications

- Text mining (more later)
- Bioinformatics
- Microarray analysis
- Mineral exploration
- Neuroscience
- Image understanding
- Air pollution research
- Weather forecasting
- ...



**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  $\mathbf{Y}$  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  $\mathbf{W}$  is fixed,  $\frac{1}{2} \|\mathbf{A} - \mathbf{WH}\|_F^2$  is convex
- Start from random  $\mathbf{W}$  and **repeat**
  - Fix  $\mathbf{W}$  and update  $\mathbf{H}$
  - Fix  $\mathbf{H}$  and update  $\mathbf{W}$
- **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:
  - $\mathbf{w}_i \leftarrow \operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\mathbf{H} - \mathbf{a}_i\|_F$
  - we can update  $\mathbf{W} \leftarrow \mathbf{A}\mathbf{H}^+$  and  $\mathbf{H} \leftarrow \mathbf{W}^+\mathbf{A}$
  - $\mathbf{X}^+$  is the pseudo-inverse of  $\mathbf{X}$  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  $\mathbf{W}$  (or  $\mathbf{H}$ ) 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  $\mathbf{W} \leftarrow [\mathbf{A}\mathbf{H}^+]_+$

# The NMF-ALS algorithm

1.  $\mathbf{W} \leftarrow \text{random}(n, k)$
2. **repeat**
  - 2.1.  $\mathbf{H} \leftarrow [\mathbf{W}^+ \mathbf{A}]_+$
  - 2.2.  $\mathbf{W} \leftarrow [\mathbf{A} \mathbf{H}^+]_+$
3. **until** convergence

# When has there been enough convergence?

- When the error doesn't change too much
  - $\| \mathbf{A} - \mathbf{W}^{(k)} \mathbf{H}^{(k)} \|_F - \| \mathbf{A} - \mathbf{W}^{(k+1)} \mathbf{H}^{(k+1)} \|_F \leq \epsilon$
- 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)
  - $f(\mathbf{H}) = \frac{1}{2} \|\mathbf{A} - \mathbf{WH}\|_F^2 = \frac{1}{2} \sum_i \|\mathbf{a}_i - \mathbf{W}\mathbf{h}_i\|_F^2$
  - $\nabla_{\mathbf{H}_{ij}} f(\mathbf{H}) = (\mathbf{W}^T \mathbf{A})_{ij} - (\mathbf{W}^T \mathbf{WH})_{ij}$
- 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.  $\mathbf{W} \leftarrow \text{random}(n, k)$
2.  $\mathbf{H} \leftarrow \text{random}(k, m)$
3. **repeat**
  - 3.1.  $\mathbf{H} \leftarrow \mathbf{H} - \varepsilon_{\mathbf{H}} \frac{\partial f}{\partial \mathbf{H}}$
  - 3.2.  $\mathbf{W} \leftarrow \mathbf{W} - \varepsilon_{\mathbf{W}} \frac{\partial f}{\partial \mathbf{W}}$
4. **until** convergence

# Oblique Projected Landweber (OPL) for NMF

- OPL provides one way to select the step size
- With  $\mathbf{H} \leftarrow \mathbf{H} - \varepsilon_{\mathbf{H}} \frac{\partial f}{\partial \mathbf{H}}$  updates, the convergence radius is  $2/\lambda_{\max}(\mathbf{W}^T \mathbf{W})$ , where  $\lambda_{\max}$  is the largest eigenvalue
  - $\lambda_{\max} \leq \max(\text{rowSums}(\mathbf{W}^T \mathbf{W}))$
- We can set the learning rates to  $1/\text{rowSums}(\mathbf{W}^T \mathbf{W})$  for a good convergence

# The OPL algorithm for updating $H$

1.  $\boldsymbol{\eta} \leftarrow \text{diag}(1 / \text{rowSums}(\mathbf{W}^T \mathbf{W}))$
2. **repeat**
  - 2.1.  $\mathbf{G} \leftarrow \mathbf{W}^T \mathbf{W} \mathbf{H} - \mathbf{W}^T \mathbf{A}$
  - 2.2.  $\mathbf{H} \leftarrow [\mathbf{H} - \boldsymbol{\eta} \mathbf{G}]_+$
3. **until** a stopping criterion is met

(small) number of iterations  
OR  
 $H$  doesn't change much

# Interior Point Gradient (IPG) for NMF

- In OPL, we might (temporarily) have negative values in  $\mathbf{W}$  or  $\mathbf{H}$
- In **Interior Point Gradient** (IPG) algorithm, we set the step sizes so that we never update to negative

# The IPG algorithm for updating $\mathbf{H}$

1. **repeat until** a stopping criterion is met

1.1.  $\mathbf{G} \leftarrow \mathbf{W}^T (\mathbf{W}\mathbf{H} - \mathbf{A})$

Gradient

1.2.  $\mathbf{D} \leftarrow \mathbf{H} / (\mathbf{W}^T \mathbf{W}\mathbf{H})$

Scaling

1.3.  $\mathbf{P} \leftarrow -\mathbf{D} * \mathbf{G}$

Update direction

1.4.  $\eta^* \leftarrow -\langle \text{vec}(\mathbf{P}), \text{vec}(\mathbf{G}) \rangle / \langle \text{vec}(\mathbf{W}\mathbf{P}), \text{vec}(\mathbf{W}\mathbf{P}) \rangle$

Best step size

1.5.  $\eta' \leftarrow \max\{\eta : \mathbf{H} + \eta\mathbf{P} \geq 0\}$

Positive step size

1.6.  $\eta \leftarrow \min\{\tau\eta', \eta^*\}$

1.7.  $\mathbf{H} \leftarrow \mathbf{H} + \eta\mathbf{P}$

Update

/ and \* are element-wise

# Multiplicative updates

- The KKT conditions for  $\mathbf{H}$  in NMF are
  - $\mathbf{H} \geq 0$ ;  $\nabla_{\mathbf{H}} \|\mathbf{A} - \mathbf{WH}\|^2/2 \geq 0$
  - $\mathbf{H} * \nabla_{\mathbf{H}} \|\mathbf{A} - \mathbf{WH}\|^2/2 = 0$  \* is element-wise product
- Substituting  $\nabla_{\mathbf{H}} \|\mathbf{A} - \mathbf{WH}\|^2/2 = \mathbf{W}^T \mathbf{WH} - \mathbf{W}^T \mathbf{A}$   
one gets  $\mathbf{H} * (\mathbf{W}^T \mathbf{WH}) = \mathbf{H} * (\mathbf{W}^T \mathbf{A})$
- This gives us an update rule for  $\mathbf{H}$

# The NMF multiplicative updates algorithm

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

2.  $\mathbf{H} \leftarrow \text{random}(k, m)$

3. **repeat**

$$3.1. \quad h_{ij} \leftarrow h_{ij} \frac{(\mathbf{W}^T \mathbf{A})_{ij}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{ij} + \varepsilon}$$

$$3.2. \quad w_{ij} \leftarrow w_{ij} \frac{(\mathbf{A} \mathbf{H}^T)_{ij}}{(\mathbf{W} \mathbf{H} \mathbf{H}^T)_{ij} + \varepsilon}$$

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

# Summary

- NMF can provide factorizations that are more interpretable than those given by SVD
- Harder to compute than SVD, but many different approaches
  - Or are they so different...
- In two weeks: Applications & alternations of NMF... *Stay tuned!*