

NONNEGATIVE MATRIX FACTORIZATION

THEORY AND APPLICATION

Presentation for Matrix Decompositions

PhD in Statistical Sciences
Academic year 2016/2017
University of Bologna

elena.geminiani4@unibo.it

30th March 2017

Table of contents

1 Nonnegative Matrix Factorization

- Introduction
- The method
- Matrix factorization
- An ill-posed problems
- Extensions

2 Application to Text Mining

- Introduction
- Analysis

3 Conclusion

Introduction

- ▶ Recent developments in technology result in increasing quantities of data, rapidly overwhelming many of classical analysis tools

Introduction

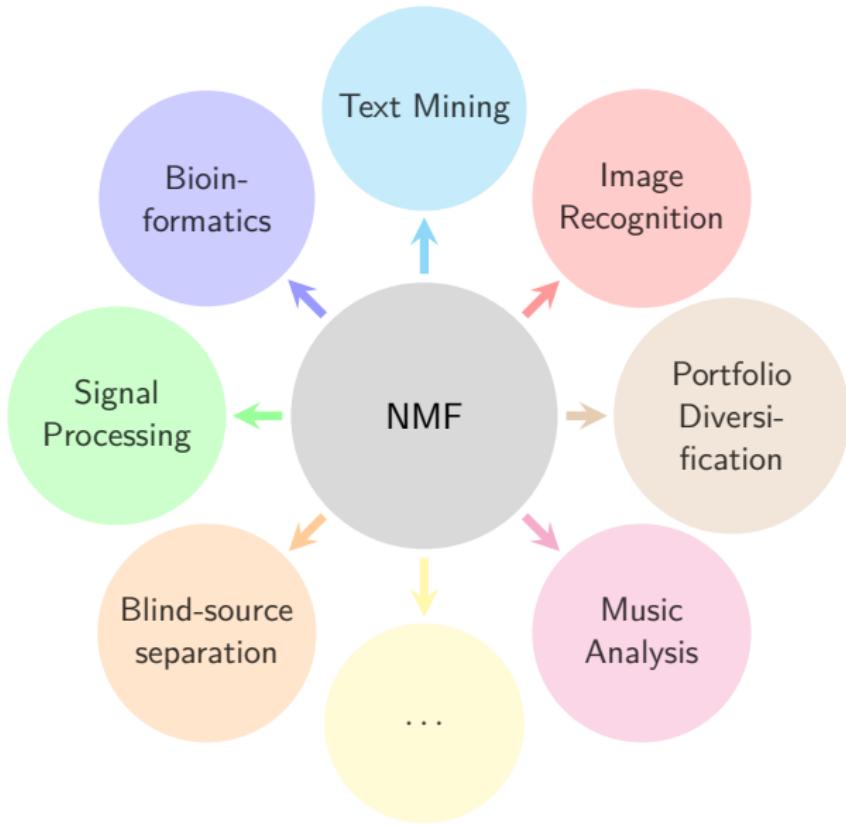
- ▶ Recent developments in technology result in increasing quantities of data, rapidly overwhelming many of classical analysis tools
 - Powerful methods for data representation and dimension reduction

Introduction

- ▶ Recent developments in technology result in increasing quantities of data, rapidly overwhelming many of classical analysis tools
 - Powerful methods for data representation and dimension reduction
- ▶ *Nonnegative data* (e.g. pixel intensities, users scores) cannot be analyzed by means of well-known techniques since they cannot guarantee to preserve nonnegativity

- ▶ Recent developments in technology result in increasing quantities of data, rapidly overwhelming many of classical analysis tools
 - Powerful methods for data representation and dimension reduction
- ▶ *Nonnegative data* (e.g. pixel intensities, users scores) cannot be analyzed by means of well-known techniques since they cannot guarantee to preserve nonnegativity
- ▶ **Nonnegative Matrix Factorization** (NMF; [Paatero and Tapper 1994](#); [Lee and Seung 1999](#)) is one of the most popular low-rank approximations obtained via subspace approximation, used for compression, visualization, feature selection and noise filtering.

Fields of application



The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

Target matrix V ($n \times m$)

$$\begin{matrix} & \text{entity 1} & \dots & \text{entity } \mu & \dots & \text{entity } m \\ \text{variable 1} & \left[\begin{array}{ccccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \text{variable } i & \left[\begin{array}{ccccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \\ \vdots & & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \text{variable } n & \left[\begin{array}{ccccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \right] \end{matrix}$$

V

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

*Target matrix V ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) **basis matrix W***

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \hline \text{variable 1} \quad \left[\begin{array}{cccc} V_{11} & \dots & V_{1\mu} & \dots & V_{1m} \end{array} \right] \quad \approx \quad \left[\begin{array}{ccccc} W_{11} & \dots & W_{1a} & \dots & W_{1r} \end{array} \right] \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \hline \text{variable } i \quad \left[\begin{array}{cccc} V_{i1} & \dots & V_{i\mu} & \dots & V_{im} \end{array} \right] \quad \approx \quad \left[\begin{array}{ccccc} W_{i1} & \dots & W_{ia} & \dots & W_{ir} \end{array} \right] \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \hline \text{variable } n \quad \left[\begin{array}{cccc} V_{n1} & \dots & V_{n\mu} & \dots & V_{nm} \end{array} \right] \quad \approx \quad \left[\begin{array}{ccccc} W_{n1} & \dots & W_{na} & \dots & W_{nr} \end{array} \right] \end{array}$$

$$V \approx W$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}*

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \quad \approx \quad \left[\begin{array}{cccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \end{array} \right] \\ \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \quad \approx \quad \left[\begin{array}{cccc} w_{i1} & \dots & w_{ia} & \dots & w_{ir} \end{array} \right] \quad \times \quad \left[\begin{array}{cccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \\ \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \quad \approx \quad \left[\begin{array}{cccc} w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array}$$

$$\mathbf{V} \approx \mathbf{WH} \quad \text{with } (\mathbf{W}, \mathbf{H}) \geq 0$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

*Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}**

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \\ \vdots \\ \vdots \\ \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \\ \vdots \\ \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{ccccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array} \times \begin{array}{c} \text{Mixture coefficients} \\ \left[\begin{array}{ccccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

$$v_{i\mu} \approx (\mathbf{WH})_{i\mu}$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

*Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}**

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \\ \vdots \\ \vdots \\ \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \\ \vdots \\ \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{ccccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array} \times \begin{array}{c} \text{Mixture coefficients} \\ \left[\begin{array}{ccccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

$$v_{i\mu} \approx (\mathbf{WH})_{i\mu} = \sum_{a=1}^r w_{ia}$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}*

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \\ \vdots \\ \vdots \\ \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \\ \vdots \\ \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{ccccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array} \times \begin{array}{c} \text{Mixture coefficients} \\ \left[\begin{array}{ccccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

$$v_{i\mu} \approx (\mathbf{WH})_{i\mu} = \sum_{a=1}^r w_{ia} h_{a\mu}$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

*Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}**

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \\ \vdots \\ \vdots \\ \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \\ \vdots \\ \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{ccccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array} \times \begin{array}{c} \text{Mixture coefficients} \\ \left[\begin{array}{ccccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

$$v_{i\mu} \approx (\mathbf{WH})_{i\mu} = \sum_{a=1}^r w_{ia} h_{a\mu}, \quad \text{for } i = 1, \dots, n, \text{ and } \mu = 1, \dots, m$$

The method

NMF allows to decompose multivariate data and search the latent features responsible for the generation of observable variables.

Target matrix \mathbf{V} ($n \times m$) is approximately factorized into the product of two nonnegative factor matrices, a ($n \times r$) *basis matrix \mathbf{W}* and a ($r \times m$) *mixture coefficients matrix \mathbf{H}*

$$\begin{array}{c} \text{entity 1} \quad \dots \quad \text{entity } \mu \quad \dots \quad \text{entity } m \\ \text{variable 1} \quad \left[\begin{array}{cccc} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{array} \right] \\ \vdots \\ \vdots \\ \vdots \\ \text{variable } i \quad \left[\begin{array}{cccc} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{array} \right] \\ \vdots \\ \vdots \\ \text{variable } n \quad \left[\begin{array}{cccc} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \end{array} \approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{ccccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \end{array} \times \begin{array}{c} \text{Mixture coefficients} \\ \left[\begin{array}{ccccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

$$v_{i\mu} \approx (\mathbf{WH})_{i\mu} = \sum_{a=1}^r w_{ia} h_{a\mu}, \quad \text{for } i = 1, \dots, n, \text{ and } \mu = 1, \dots, m$$

A crucial role is played by the factorization rank $r \ll \min(n, m)$ such that $(n+m)r < nm$.

Matrix Factorization

The factorization of the matrix V can be achieved in several ways using numerical algorithms (block coordinate descendent scheme).

Matrix Factorization

The factorization of the matrix \mathbf{V} can be achieved in several ways using numerical algorithms (block coordinate descendent scheme).

General framework of NMF algorithms

Input: A nonnegative matrix $\mathbf{V} \in \mathbb{R}_+^{n \times m}$ and the factorization rank r

- ▶ Generate some initial matrices $\mathbf{W}^{(0)} \geq 0, \mathbf{H}^{(0)} \geq 0$
- ▶ **for** $t = 1, 2, \dots$, stopping time
 - $\mathbf{W}^{(t)} = \text{update } (\mathbf{V}, \mathbf{H}^{(t-1)}, \mathbf{W}^{(t-1)})$
 - $\mathbf{H}^{(t)^T} = \text{update } (\mathbf{V}^T, \mathbf{W}^{(t)^T}, \mathbf{H}^{(t-1)^T})$
- end for

Output: A rank- r NMF of $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, for $(\mathbf{W}, \mathbf{H}) \geq 0$

Matrix Factorization

The factorization of the matrix \mathbf{V} can be achieved in several ways using numerical algorithms (block coordinate descendent scheme).

General framework of NMF algorithms

Input: A nonnegative matrix $\mathbf{V} \in \mathbb{R}_+^{n \times m}$ and the factorization rank r

- ▶ Generate some initial matrices $\mathbf{W}^{(0)} \geq 0, \mathbf{H}^{(0)} \geq 0$
- ▶ **for** $t = 1, 2, \dots$, stopping time
 - $\mathbf{W}^{(t)} = \text{update } (\mathbf{V}, \mathbf{H}^{(t-1)}, \mathbf{W}^{(t-1)})$
 - $\mathbf{H}^{(t)^T} = \text{update } (\mathbf{V}^T, \mathbf{W}^{(t)^T}, \mathbf{H}^{(t-1)^T})$
- end for

Output: A rank- r NMF of $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, for $(\mathbf{W}, \mathbf{H}) \geq 0$

- ➡ Objective function
- ➡ Numerical algorithm

Objective function

To find an approximate factorization $\mathbf{V} \approx \mathbf{WH}$, one needs to define loss functions $L(\mathbf{V}, \mathbf{WH})$ that quantify the quality of the approximation.

Objective function

To find an approximate factorization $\mathbf{V} \approx \mathbf{WH}$, one needs to define loss functions $L(\mathbf{V}, \mathbf{WH})$ that quantify the quality of the approximation.

- ▶ Square Euclidean distance:

$$L(\mathbf{V}, \mathbf{WH}) = \|\mathbf{V} - \mathbf{WH}\|^2 = \sum_{i=1}^n \sum_{\mu=1}^m (v_{i\mu} - (\mathbf{WH})_{i\mu})^2$$

Objective function

To find an approximate factorization $\mathbf{V} \approx \mathbf{WH}$, one needs to define loss functions $L(\mathbf{V}, \mathbf{WH})$ that quantify the quality of the approximation.

- ▶ Square Euclidean distance:

$$L(\mathbf{V}, \mathbf{WH}) = \|\mathbf{V} - \mathbf{WH}\|^2 = \sum_{i=1}^n \sum_{\mu=1}^m (v_{i\mu} - (\mathbf{WH})_{i\mu})^2$$

- ▶ Generalized Kullback-Leibler divergence

$$L(\mathbf{V}, \mathbf{WH}) = D(\mathbf{V} || \mathbf{WH}) = \sum_{i=1}^n \sum_{\mu=1}^m v_{i\mu} \log \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}} - v_{i\mu} + (\mathbf{WH})_{i\mu}$$

Objective function

To find an approximate factorization $\mathbf{V} \approx \mathbf{WH}$, one needs to define loss functions $L(\mathbf{V}, \mathbf{WH})$ that quantify the quality of the approximation.

- Square Euclidean distance:

$$L(\mathbf{V}, \mathbf{WH}) = \|\mathbf{V} - \mathbf{WH}\|^2 = \sum_{i=1}^n \sum_{\mu=1}^m (v_{i\mu} - (\mathbf{WH})_{i\mu})^2$$

- Generalized Kullback-Leibler divergence

$$L(\mathbf{V}, \mathbf{WH}) = D(\mathbf{V} || \mathbf{WH}) = \sum_{i=1}^n \sum_{\mu=1}^m v_{i\mu} \log \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}} - v_{i\mu} + (\mathbf{WH})_{i\mu}$$

Nonlinear optimization problem: minimize the objective function with respect to matrix factors, subject to nonnegativity constraint:

$$\min_{\mathbf{W}, \mathbf{H} \geq 0} \left[\underbrace{L(\mathbf{V}, \mathbf{WH})}_{\text{Loss function}} + \underbrace{R(\mathbf{V}, \mathbf{WH})}_{\text{Regularization term}} \right]$$

Matrix Factorization

The factorization of the matrix \mathbf{V} can be achieved in several ways using numerical algorithms (block coordinate descendent scheme).

General framework of NMF algorithms

Input: A nonnegative matrix $\mathbf{V} \in \mathbb{R}_+^{n \times m}$ and the factorization rank r

- ▶ Generate some initial matrices $\mathbf{W}^{(0)} \geq 0, \mathbf{H}^{(0)} \geq 0$
- ▶ **for** $t = 1, 2, \dots$, stopping time
 - $\mathbf{W}^{(t)} = \text{update } (\mathbf{V}, \mathbf{H}^{(t-1)}, \mathbf{W}^{(t-1)})$
 - $\mathbf{H}^{(t)^T} = \text{update } (\mathbf{V}^T, \mathbf{W}^{(t)^T}, \mathbf{H}^{(t-1)^T})$
- end for

Output: A rank- r NMF of $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, for $(\mathbf{W}, \mathbf{H}) \geq 0$

Matrix Factorization

The factorization of the matrix \mathbf{V} can be achieved in several ways using numerical algorithms (block coordinate descendent scheme).

General framework of NMF algorithms

Input: A nonnegative matrix $\mathbf{V} \in \mathbb{R}_+^{n \times m}$ and the factorization rank r

- ▶ Generate some initial matrices $\mathbf{W}^{(0)} \geq 0, \mathbf{H}^{(0)} \geq 0$
- ▶ **for** $t = 1, 2, \dots$, stopping time
 - $\mathbf{W}^{(t)} = \text{update } (\mathbf{V}, \mathbf{H}^{(t-1)}, \mathbf{W}^{(t-1)})$
 - $\mathbf{H}^{(t)^T} = \text{update } (\mathbf{V}^T, \mathbf{W}^{(t)^T}, \mathbf{H}^{(t-1)^T})$
- end for

Output: A rank- r NMF of $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, for $(\mathbf{W}, \mathbf{H}) \geq 0$

➡ Objective function

➡ Numerical algorithm

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H}

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and $\mathbf{H} \rightarrow$ local minima

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and $\mathbf{H} \rightarrow$ local minima

Numerical techniques can be classified in three categories:

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and $\mathbf{H} \rightarrow$ local minima

Numerical techniques can be classified in three categories:

- ① **Gradient descent algorithms:** simplest to implement, but very slow to converge (if at all); sensitive to the step size choice and lack a convergence supporting theory;

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H} → local minima

Numerical techniques can be classified in three categories:

- ① **Gradient descent algorithms:** simplest to implement, but very slow to converge (if at all); sensitive to the step size choice and lack a convergence supporting theory;
- ② **Multiplicative update rules:** extensively used due to their simplicity, applicable to sparse matrices, proposed as numerical algorithm for solving NMF by [Lee and Seung 2001](#);

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H} → local minima

Numerical techniques can be classified in three categories:

- ① **Gradient descent algorithms**: simplest to implement, but very slow to converge (if at all); sensitive to the step size choice and lack a convergence supporting theory;
- ② **Multiplicative update rules**: extensively used due to their simplicity, applicable to sparse matrices, proposed as numerical algorithm for solving NMF by [Lee and Seung 2001](#);
- ③ **Alternating least squares algorithms** ([Kim and Park 2007](#)): a least squares step is followed by another least squares step in an alternating fashion; very fast, quick to converge, and give accurate factors.

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H} → local minima

Numerical techniques can be classified in three categories:

- ✗ **Gradient descent algorithms:** simplest to implement, but very slow to converge (if at all); sensitive to the step size choice and lack a convergence supporting theory;
- ✓ **Multiplicative update rules:** extensively used due to their simplicity, applicable to sparse matrices, proposed as numerical algorithm for solving NMF by [Lee and Seung 2001](#);
- ✓ **Alternating least squares algorithms** ([Kim and Park 2007](#)): a least squares step is followed by another least squares step in an alternating fashion; very fast, quick to converge, and give accurate factors.

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H} → local minima

Numerical techniques can be classified in three categories:

- ① **Gradient descent algorithms**: simplest to implement, but very slow to converge (if at all); sensitive to the step size choice and lack a convergence supporting theory;
- **Multiplicative update rules**: extensively used due to their simplicity, applicable to sparse matrices, proposed as numerical algorithm for solving NMF by [Lee and Seung 2001](#);
- ② **Alternating least squares algorithms** ([Kim and Park 2007](#)): a least squares step is followed by another least squares step in an alternating fashion; very fast, quick to converge, and give accurate factors.

Multiplicative rules

They can be divided according to the choice of the objective function.

For $a = 1, \dots, r$ $\mu = 1, \dots, m$ $i = 1, \dots, n$:

Multiplicative rules

They can be divided according to the choice of the objective function.

For $a = 1, \dots, r$ $\mu = 1, \dots, m$ $i = 1, \dots, n$:

Divergence measure

$$h_{a\mu} \leftarrow h_{a\mu} \frac{\sum_{i=1}^n w_{ia} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{i=1}^n w_{ia}}$$

$$w_{ia} \leftarrow w_{ia} \frac{\sum_{\mu=1}^m h_{a\mu} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{\mu=1}^m h_{a\mu}}$$

Multiplicative rules

They can be divided according to the choice of the objective function.

For $a = 1, \dots, r$ $\mu = 1, \dots, m$ $i = 1, \dots, n$:

Divergence measure

$$h_{a\mu} \leftarrow h_{a\mu} \frac{\sum_{i=1}^n w_{ia} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{i=1}^n w_{ia}}$$
$$w_{ia} \leftarrow w_{ia} \frac{\sum_{\mu=1}^m h_{a\mu} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{\mu=1}^m h_{a\mu}}$$

Square Euclidean distance

$$h_{a\mu} \leftarrow h_{a\mu} \frac{(\mathbf{W}^T \mathbf{V})_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$$
$$w_{ia} \leftarrow w_{ia} \frac{(\mathbf{V} \mathbf{H}^T)_{ia}}{(\mathbf{W} \mathbf{H} \mathbf{H}^T)_{ia}}$$

Multiplicative rules

They can be divided according to the choice of the objective function.

For $a = 1, \dots, r$ $\mu = 1, \dots, m$ $i = 1, \dots, n$:

Divergence measure

$$h_{a\mu} \leftarrow h_{a\mu} \frac{\sum_{i=1}^n w_{ia} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{i=1}^n w_{ia}}$$
$$w_{ia} \leftarrow w_{ia} \frac{\sum_{\mu=1}^m h_{a\mu} \frac{v_{i\mu}}{(\mathbf{WH})_{i\mu}}}{\sum_{\mu=1}^m h_{a\mu}}$$

Square Euclidean distance

$$h_{a\mu} \leftarrow h_{a\mu} \frac{(\mathbf{W}^T \mathbf{V})_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$$
$$w_{ia} \leftarrow w_{ia} \frac{(\mathbf{V} \mathbf{H}^T)_{ia}}{(\mathbf{W} \mathbf{H} \mathbf{H}^T)_{ia}}$$

Objective function guaranteed to monotonically decrease, but not necessarily converges to a local minimum.

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and $\mathbf{H} \rightarrow$ local minima

Numerical techniques can be classified in three categories:

Numerical algorithms

Cost functions are not simultaneously convex in both \mathbf{W} and \mathbf{H} → local minima

Numerical techniques can be classified in three categories:

- ① **Gradient descent algorithms**: simplest to implement, but very slow to converge (if at all). Their convergence is very sensitive to the step size choice and a convergence supporting theory is missing;
- ② **Multiplicative update rules**: extensively used since they are simple to implement, applicable to sparse matrices, and they were proposed as numerical algorithm for solving NMF by Lee and Seung 2001;
- **Alternating least squares algorithms** (Kim and Park 2007) : a least squares step is followed by another least squares step in an alternating fashion; very fast, quick to convergence, and give accurate factors.

Alternating least squares

Developed to solve the issues affecting multiplicative rules:

- ▶ Slow convergence → alternating strategy

Alternating least squares

Developed to solve the issues affecting multiplicative rules:

- ▶ Slow convergence → alternating strategy
- ▶ Locking property → escaping from poor local minima

Alternating least squares

Developed to solve the issues affecting multiplicative rules:

- ▶ Slow convergence → alternating strategy
- ▶ Locking property → escaping from poor local minima

Basic ALS

- ① Take \mathbf{W} as fixed
- ② Solve for \mathbf{H} in the equation $\mathbf{W}^T \mathbf{W} \mathbf{H} = \mathbf{W}^T \mathbf{V}$
- ③ Set all negative elements in \mathbf{H} to zero
- ④ Then, take \mathbf{H} as fixed
- ⑤ Solves for \mathbf{W} in the equation $\mathbf{H} \mathbf{H}^T \mathbf{W}^T = \mathbf{H} \mathbf{V}^T$
- ⑥ Set all the negative elements in \mathbf{W} to zero.

Alternating least squares

Developed to solve the issues affecting multiplicative rules:

- ▶ Slow convergence → alternating strategy
- ▶ Locking property → escaping from poor local minima

Basic ALS

- ① Take \mathbf{W} as fixed
- ② Solve for \mathbf{H} in the equation $\mathbf{W}^T \mathbf{W} \mathbf{H} = \mathbf{W}^T \mathbf{V}$
- ③ Set all negative elements in \mathbf{H} to zero
- ④ Then, take \mathbf{H} as fixed
- ⑤ Solves for \mathbf{W} in the equation $\mathbf{H} \mathbf{H}^T \mathbf{W}^T = \mathbf{H} \mathbf{V}^T$
- ⑥ Set all the negative elements in \mathbf{W} to zero.

Computationally demanding due to additional nonnegativity constraint on least squares.

Initialization

NMF algorithms need to be initialized with values for $\mathbf{W}^{(0)}$ and/or $\mathbf{H}^{(0)}$.

Initialization

NMF algorithms need to be initialized with values for $\mathbf{W}^{(0)}$ and/or $\mathbf{H}^{(0)}$.

Standard NMF relies on **random initialization** → very simple, but requires multiple runs with different starting points (high computation times).

Initialization

NMF algorithms need to be initialized with values for $\mathbf{W}^{(0)}$ and/or $\mathbf{H}^{(0)}$.

Standard NMF relies on **random initialization** → very simple, but requires multiple runs with different starting points (high computation times).

Several seeding methods have been proposed, based on:

- **Independent component analysis**
- **Nonnegative double singular value decomposition**
- **Clustering techniques**

In practice, one may use several different initializations via a Monte Carlo approach, and keep the best solution obtained.

Factorization rank

- ↳ Too high values → potential risk of overfitting
- ↳ Too low values → potential risk of misinterpretation

Factorization rank

- ↳ Too high values → potential risk of overfitting
- ↳ Too low values → potential risk of misinterpretation

Consider the following application to image recognition:

Factorization rank

- ➡ Too high values → potential risk of overfitting
- ➡ Too low values → potential risk of misinterpretation

Consider the following application to image recognition:



Figure 1: *Madonna del Cardellino* in grey-scale format; Target matrix rank = 400

Factorization rank

Approximated paintings with different factorization ranks: **5**, 20, 40, 80.

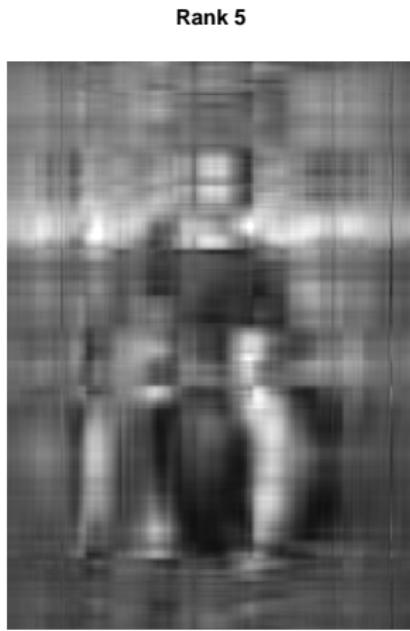


Figure 2: Factorization rank=5

Factorization rank

Approximated paintings with different factorization ranks: 5, **20**, 40, 80.

Rank 20



Figure 3: Factorization rank=20

Factorization rank

Approximated paintings with different factorization ranks: 5, 20, **40**, 80.



Figure 4: Factorization rank=40

Factorization rank

Approximated paintings with different factorization ranks: 5, 20, 40, **80**.



Figure 5: Factorization rank=80

Factorization rank

Unless one has prior knowledge on domain theory, a common way to choose r is through **trial and error**.

Factorization rank

Unless one has prior knowledge on domain theory, a common way to choose r is through **trial and error**.

However, several approaches have been proposed:

- ▶ Apply **SVD** and look at the decay of singular values of the input data matrix;

Factorization rank

Unless one has prior knowledge on domain theory, a common way to choose r is through **trial and error**.

However, several approaches have been proposed:

- ▶ Apply **SVD** and look at the decay of singular values of the input data matrix;
- ▶ Take the first value of r for which the **cophenetic coefficient** starts decreasing;

Factorization rank

Unless one has prior knowledge on domain theory, a common way to choose r is through **trial and error**.

However, several approaches have been proposed:

- ▶ Apply **SVD** and look at the decay of singular values of the input data matrix;
- ▶ Take the first value of r for which the **cophenetic coefficient** starts decreasing;
- ▶ Choose the first value where the **RSS curve** presents an inflection point;

Factorization rank

Unless one has prior knowledge on domain theory, a common way to choose r is through **trial and error**.

However, several approaches have been proposed:

- ▶ Apply **SVD** and look at the decay of singular values of the input data matrix;
- ▶ Take the first value of r for which the **cophenetic coefficient** starts decreasing;
- ▶ Choose the first value where the **RSS curve** presents an inflection point;
- ▶ Consider smallest value at which the decrease in the **RSS** is lower than the decrease of analogous curve from **random data**.

A unique decomposition?

- ❖ Find a simplicial convex cone $C_{\mathbf{W}} = \left\{ \sum_{a=1}^r \lambda_a \mathbf{w}_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

A unique decomposition?

- ❖ Find a simplicial convex cone $C_W = \left\{ \sum_{a=1}^r \lambda_a \mathbf{w}_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

If data values are strictly positive:

$$v_{i\mu} \geq \epsilon \quad \forall \epsilon > 0$$

the column vectors of \mathbf{V} are well inside the interior of positive orthant of \mathbb{R}^n , and there are many simplicial cones containing the data (Donoho and Stodden 2004).

A unique decomposition?

- ❖ Find a simplicial convex cone $C_W = \left\{ \sum_{a=1}^r \lambda_a w_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

If data values are strictly positive:

$$v_{i\mu} \geq \epsilon \quad \forall \epsilon > 0$$

the column vectors of \mathbf{V} are well inside the interior of positive orthant of \mathbb{R}^n , and there are many simplicial cones containing the data (Donoho and Stodden 2004).

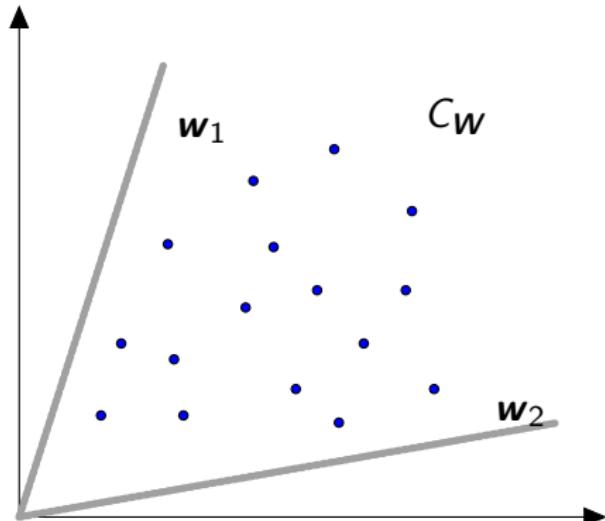


Figure 6: Non-uniqueness of nonnegative matrix factorization

A unique decomposition?

- ❖ Find a simplicial convex cone $C_W = \left\{ \sum_{a=1}^r \lambda_a w_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

If data values are strictly positive:

$$v_{i\mu} \geq \epsilon \quad \forall \epsilon > 0$$

the column vectors of \mathbf{V} are well inside the interior of positive orthant of \mathbb{R}^n , and there are many simplicial cones containing the data (Donoho and Stodden 2004).

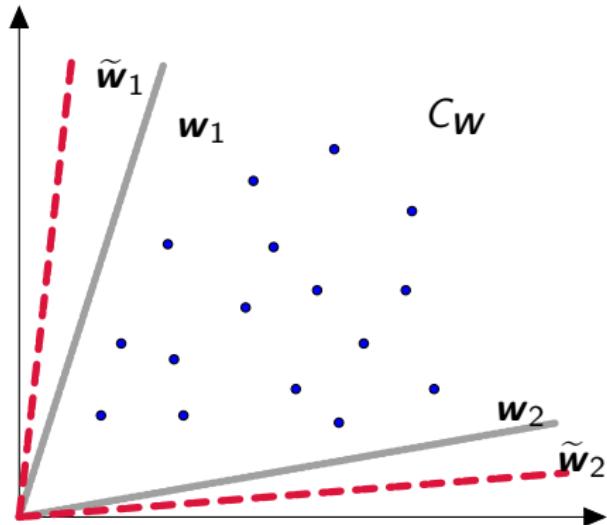


Figure 6: Non-uniqueness of nonnegative matrix factorization

A unique decomposition?

- ❖ Find a simplicial convex cone $C_{\mathbf{W}} = \left\{ \sum_{a=1}^r \lambda_a \mathbf{w}_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

Under such strict positivity condition, there are many distinct factorizations.

A unique decomposition?

- ❖ Find a simplicial convex cone $C_{\mathbf{W}} = \left\{ \sum_{a=1}^r \lambda_a \mathbf{w}_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

Under such strict positivity condition, there are many distinct factorizations.

NMF is an ill-posed problem

Any nonnegative invertible matrix \mathbf{Q} satisfying $\begin{cases} \widetilde{\mathbf{W}} = \mathbf{W}\mathbf{Q} \geq 0 \\ \widetilde{\mathbf{H}} = \mathbf{Q}^{-1}\mathbf{H} \geq 0 \end{cases}$ generates an equivalent factorization:

$$\widetilde{\mathbf{W}}\widetilde{\mathbf{H}} = \mathbf{W}\mathbf{H}$$

A unique decomposition?

- ❖ Find a simplicial convex cone $C_{\mathbf{W}} = \left\{ \sum_{a=1}^r \lambda_a \mathbf{w}_a, \lambda_a \geq 0 \right\}$ contained in the positive orthant and containing data points.

Under such strict positivity condition, there are many distinct factorizations.

NMF is an ill-posed problem

Any nonnegative invertible matrix \mathbf{Q} satisfying $\begin{cases} \widetilde{\mathbf{W}} = \mathbf{W}\mathbf{Q} \geq 0 \\ \widetilde{\mathbf{H}} = \mathbf{Q}^{-1}\mathbf{H} \geq 0 \end{cases}$ generates an equivalent factorization:

$$\widetilde{\mathbf{W}}\widetilde{\mathbf{H}} = \mathbf{W}\mathbf{H}$$

- Impose additional constraints to confine the feasible solution set.

Extensions

Over the years, the standard NMF model has been extended to explicitly include auxiliary constraints (incorporated in the objective function) on \mathbf{W} and/or \mathbf{H} to improve the factorization.

Extensions

Over the years, the standard NMF model has been extended to explicitly include auxiliary constraints (incorporated in the objective function) on \mathbf{W} and/or \mathbf{H} to improve the factorization.

- ▶ **Smoothness:** regularizes computed solutions in presence of noise in data;

Extensions

Over the years, the standard NMF model has been extended to explicitly include auxiliary constraints (incorporated in the objective function) on \mathbf{W} and/or \mathbf{H} to improve the factorization.

- ▶ **Smoothness:** regularizes computed solutions in presence of noise in data;
- ▶ **Sparseness:** provides a clearer representation of data. Sparseness is defined as ([Hoyer 2004](#)):

$$\text{sparseness}(\mathbf{x}) = \frac{\sqrt{n} - \frac{\|\mathbf{x}\|_1}{\|\mathbf{x}\|_2}}{\sqrt{n} - 1}$$

Extensions

Over the years, the standard NMF model has been extended to explicitly include auxiliary constraints (incorporated in the objective function) on \mathbf{W} and/or \mathbf{H} to improve the factorization.

- ▶ **Smoothness:** regularizes computed solutions in presence of noise in data;
- ▶ **Sparseness:** provides a clearer representation of data. Sparseness is defined as ([Hoyer 2004](#)):

$$\text{sparseness}(\mathbf{x}) = \frac{\sqrt{n} - \frac{\|\mathbf{x}\|_1}{\|\mathbf{x}\|_2}}{\sqrt{n} - 1}$$

Some other extensions of standard NMF include:

- ① *Non-smooth NMF (nsNMF);*
- ② *NMF with offset;*
- ③ *Pattern-expression NMF*

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Term-document matrix

	Doc 1	Doc μ	...	Doc m
Word 1	v_{11}	$v_{1\mu}$	v_{1m}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word i	v_{i1}	$v_{i\mu}$	v_{im}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word n	v_{n1}	$v_{n\mu}$	v_{nm}

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Term-document matrix

	Doc 1	Doc μ	...	Doc m
Word 1	v_{11}	$v_{1\mu}$	v_{1m}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word i	v_{i1}	$v_{i\mu}$	v_{im}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word n	v_{n1}	$v_{n\mu}$	v_{nm}

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Term-document matrix

$$\begin{array}{c} \text{Term-document matrix} \\ \text{Basis vectors} \\ \approx \end{array}$$
$$\begin{matrix} & \text{Doc 1} & \dots & \text{Doc } \mu & \dots & \text{Doc } m \\ \text{Word 1} & \left[\begin{matrix} v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \end{matrix} \right] & \approx & \left[\begin{matrix} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \text{Word } i & \left[\begin{matrix} v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \end{matrix} \right] & \approx & \left[\begin{matrix} w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \text{Word } n & \left[\begin{matrix} v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{matrix} \right] & \approx & \left[\begin{matrix} w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{matrix} \right] \end{matrix} \right. \end{matrix}$$

The diagram illustrates the decomposition of a term-document matrix into basis vectors. On the left, the term-document matrix is shown as a grid of words (rows) by documents (columns). The entry $v_{i\mu}$ is highlighted with a blue box. The matrix is approximately equal (\approx) to a product of two matrices on the right: the basis vectors matrix and the document-topic matrix. The basis vectors matrix has columns labeled $w_{11}, \dots, w_{1a}, \dots, w_{1r}$, $w_{i1}, \dots, w_{ia}, \dots, w_{ir}$, and $w_{n1}, \dots, w_{na}, \dots, w_{nr}$. The document-topic matrix has columns labeled $w_{11}, \dots, w_{1a}, \dots, w_{1r}$, $w_{i1}, \dots, w_{ia}, \dots, w_{ir}$, and $w_{n1}, \dots, w_{na}, \dots, w_{nr}$.

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Term-document matrix

$$\begin{array}{c} \text{Term-document matrix} \\ \text{Basis vectors} \\ \approx \\ \begin{array}{l} \text{Word 1} \\ \vdots \\ \text{Word } i \\ \vdots \\ \text{Word } n \end{array} \left[\begin{array}{cccc} \text{Doc 1} & \dots & \text{Doc } \mu & \dots & \text{Doc } m \\ v_{11} & \dots & v_{1\mu} & \dots & v_{1m} \\ \vdots & & \vdots & & \vdots \\ v_{i1} & \dots & v_{i\mu} & \dots & v_{im} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ v_{n1} & \dots & v_{n\mu} & \dots & v_{nm} \end{array} \right] \right] \approx \left[\begin{array}{cccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \\ \uparrow \quad \uparrow \quad \uparrow \\ \text{Latent Topics} \end{array}$$

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.

Term-document matrix

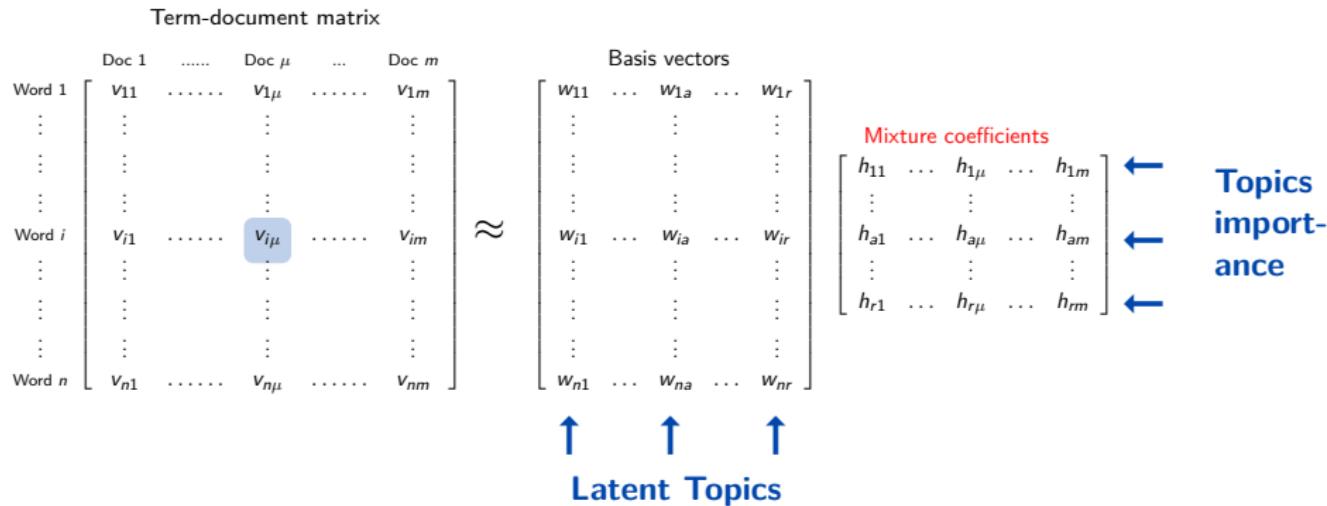
	Doc 1	Doc μ	...	Doc m
Word 1	v_{11}	$v_{1\mu}$	v_{1m}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word i	v_{i1}	$v_{i\mu}$	v_{im}
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
⋮	⋮		⋮		⋮
Word n	v_{n1}	$v_{n\mu}$	v_{nm}

$$\approx \begin{array}{c} \text{Basis vectors} \\ \left[\begin{array}{cccc} w_{11} & \dots & w_{1a} & \dots & w_{1r} \\ \vdots & & \vdots & & \vdots \\ w_{i1} & \dots & w_{ia} & \dots & w_{ir} \\ \vdots & & \vdots & & \vdots \\ w_{n1} & \dots & w_{na} & \dots & w_{nr} \end{array} \right] \\ \text{Mixture coefficients} \\ \left[\begin{array}{cccc} h_{11} & \dots & h_{1\mu} & \dots & h_{1m} \\ \vdots & & \vdots & & \vdots \\ h_{a1} & \dots & h_{a\mu} & \dots & h_{am} \\ \vdots & & \vdots & & \vdots \\ h_{r1} & \dots & h_{r\mu} & \dots & h_{rm} \end{array} \right] \end{array}$$

↑ ↑ ↑
Latent Topics

Application to Text Mining

Text analysis learns from collections of text data to recover topics and classify documents.



Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);

Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;

Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;
- ③ Term-document matrix of ≈ 19000 words; sparseness degree $\approx 95\%$.

Analysis

- ① Analyzed texts: complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;
- ③ Term-document matrix of ≈ 19000 words; sparseness degree $\approx 95\%$.



Figure 7: Word cloud (words occurring more than 1000 times)

Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;
- ③ Term-document matrix of ≈ 19000 words; sparseness degree $\approx 95\%$.
- ④ Removal of unfrequent terms → 63×182 term-document matrix;

Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;
- ③ Term-document matrix of ≈ 19000 words; sparseness degree $\approx 95\%$.
- ④ Removal of unfrequent terms → 63×182 term-document matrix;
- ⑤ A rank-3 NMF model fitted to the data using multiplicative update minimizing Kullback-Leibler divergence (model estimated 200 times, each time starting from different set of initial values);

Analysis

- ① *Analyzed texts:* complete published works by William Shakespeare (corpus composed by 182 docs);
- ② Variety of preprocessing operations;
- ③ Term-document matrix of ≈ 19000 words; sparseness degree $\approx 95\%$.
- ④ Removal of unfrequent terms → 63×182 term-document matrix;
- ⑤ A rank-3 NMF model fitted to the data using multiplicative update minimizing Kullback-Leibler divergence (model estimated 200 times, each time starting from different set of initial values);
- ⑥ Factor matrices explored by means of heatmaps.

Topic recovery

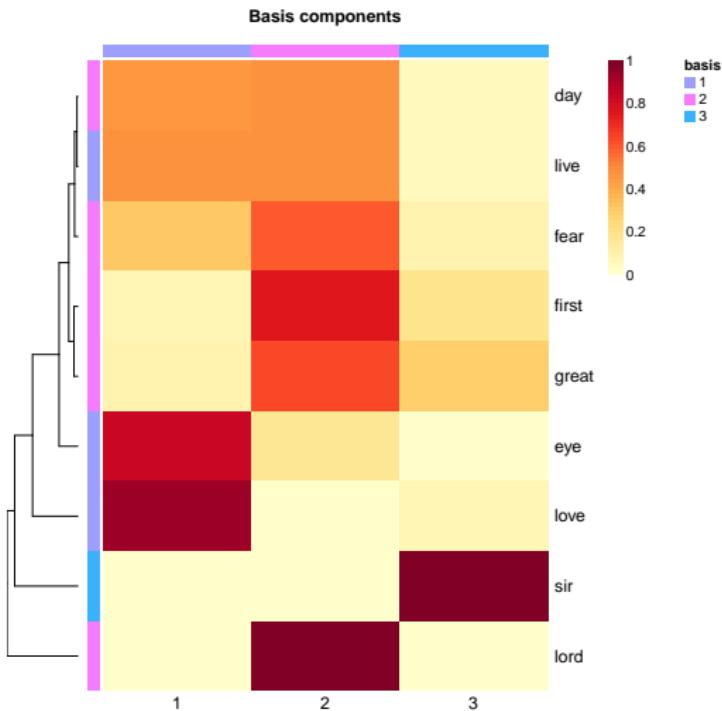
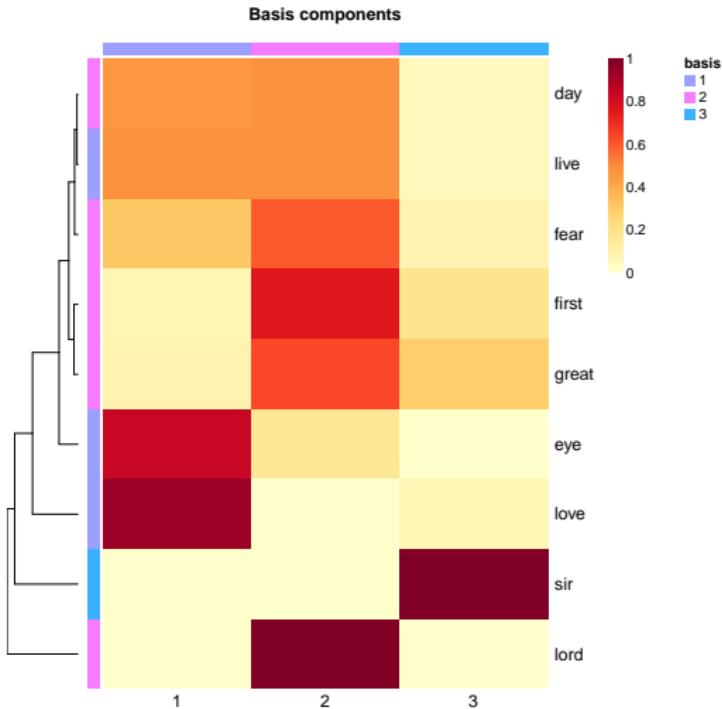


Figure 7: Heatmap of basis matrix (only relevant rows selected)

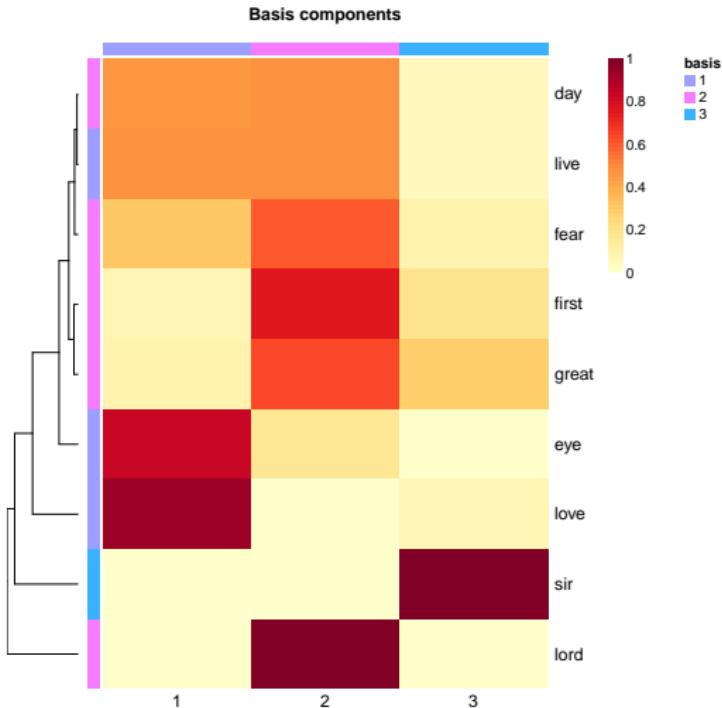
Topic recovery



- First basis vector: entails to common topic of **love, beauty, and eternity**;

Figure 7: Heatmap of basis matrix (only relevant rows selected)

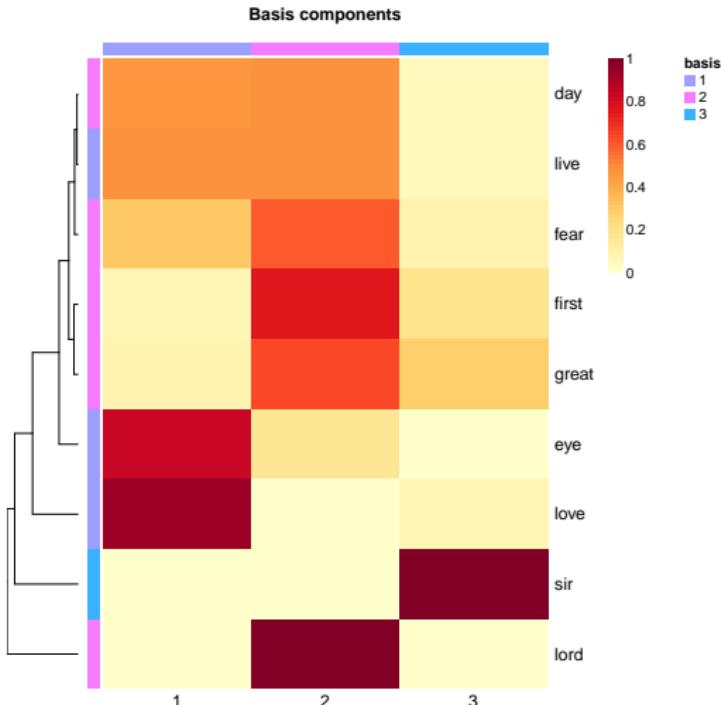
Topic recovery



- First basis vector: entails to common topic of **love, beauty, and eternity**;
- Second basis vector: activated by words standing for various themes: **obeisance, grandeur**, but also **power** and **impending death**;

Figure 7: Heatmap of basis matrix (only relevant rows selected)

Topic recovery



- First basis vector: entails to common topic of **love, beauty, and eternity**;
- Second basis vector: activated by words standing for various themes: **obeisance, grandeur**, but also **power** and **impending death**;
- Third basis component: related to comicality, **sarcasm** and **irony**.

Figure 7: Heatmap of basis matrix (only relevant rows selected)

Document classification

Allocate docs to basis component for which they have the highest coefficient.

Document classification

Allocate docs to basis component for which they have the highest coefficient.

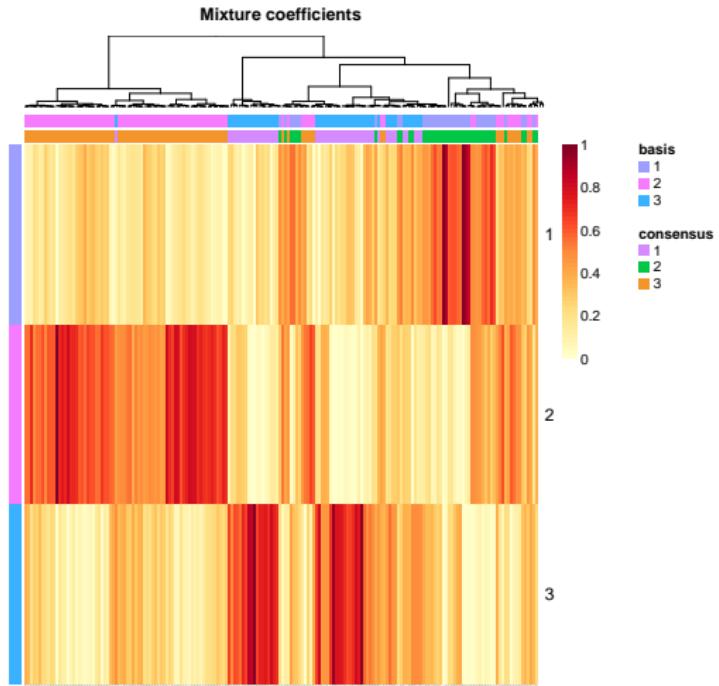


Figure 8: Mixture coefficients heatmap

Document classification

Allocate docs to basis component for which they have the highest coefficient.

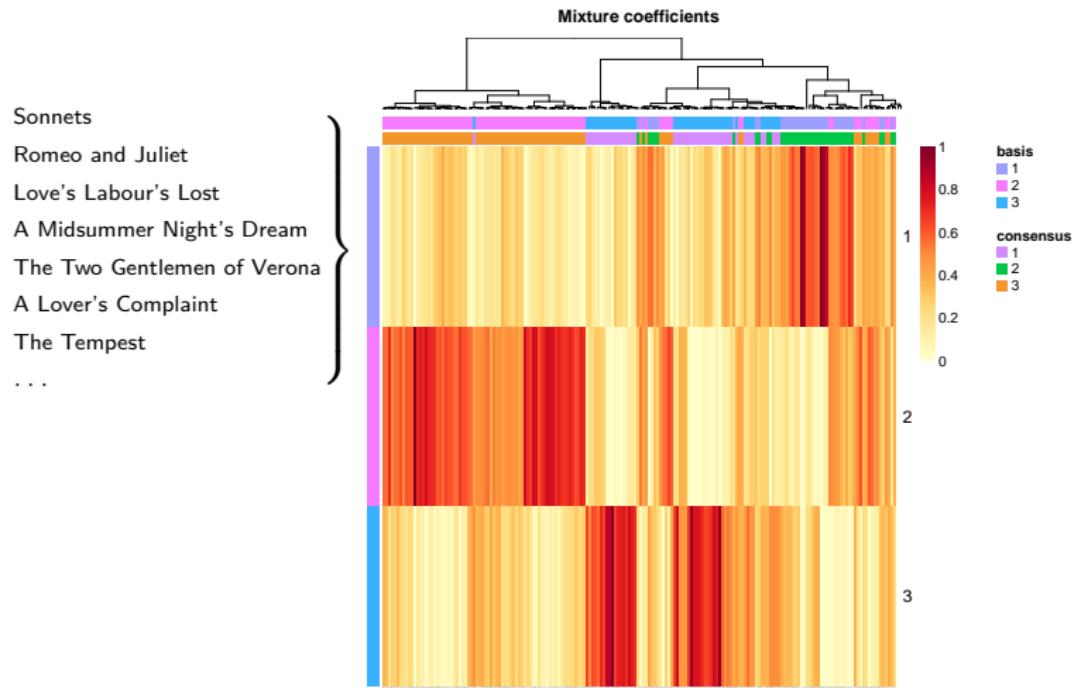


Figure 8: Mixture coefficients heatmap

Document classification

Allocate docs to basis component for which they have the highest coefficient.

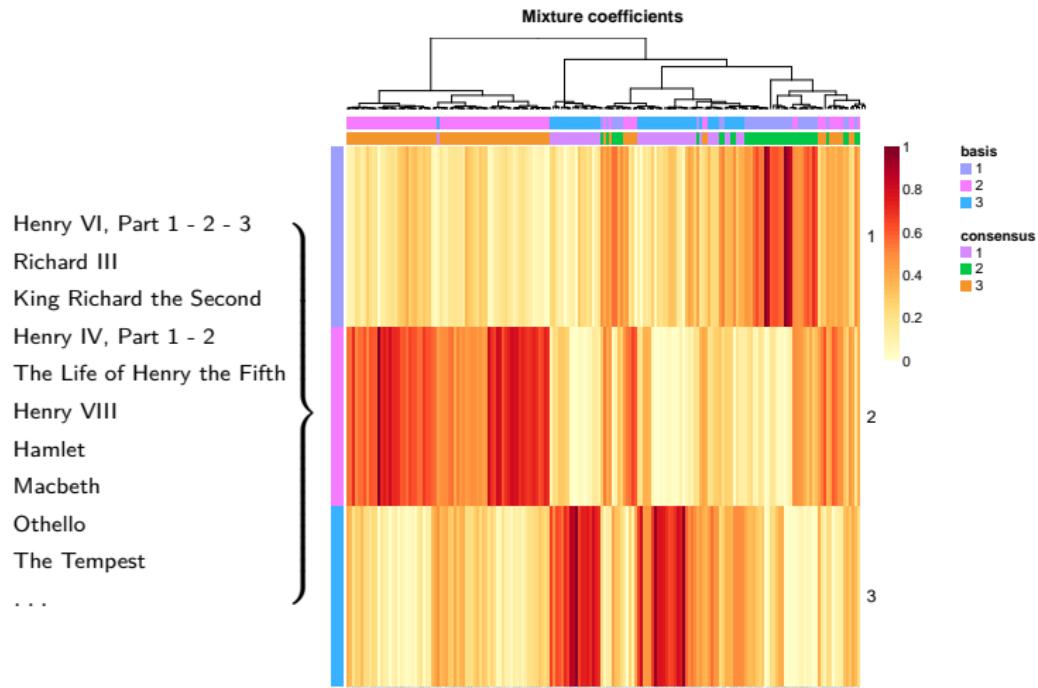


Figure 8: Mixture coefficients heatmap

Document classification

Allocate docs to basis component for which they have the highest coefficient.

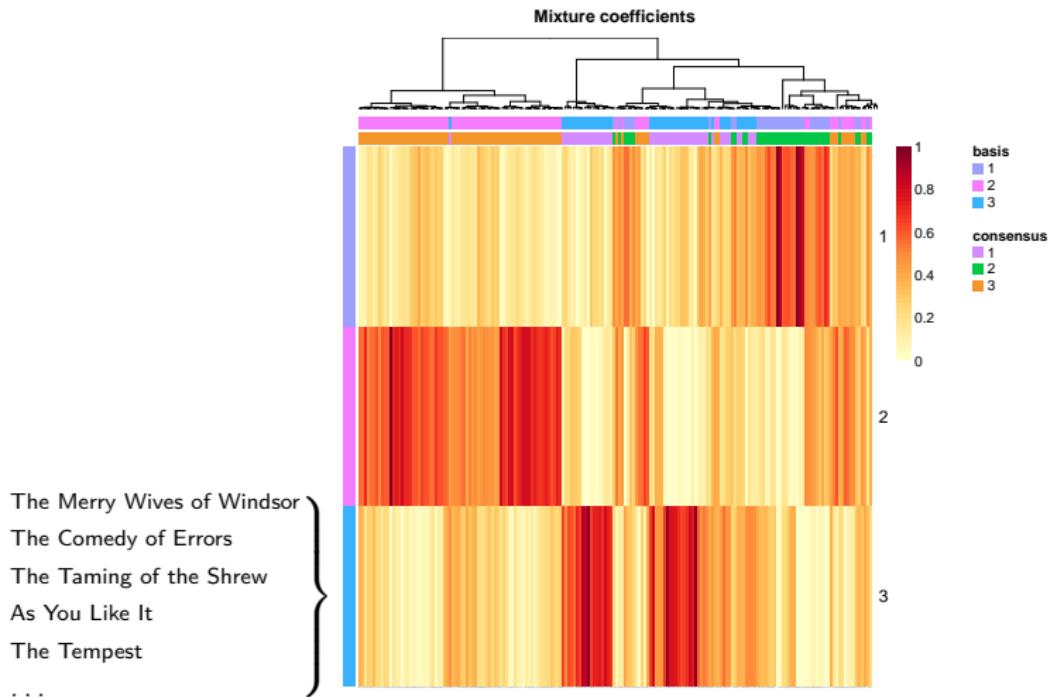


Figure 8: Mixture coefficients heatmap

Conclusion

NMF has made great strides and huge impact in the last few years and, by looking at the number of yearly citations received by founding paper of [Lee and Seung 1999](#), it's very likely it will have a promising and bright future.

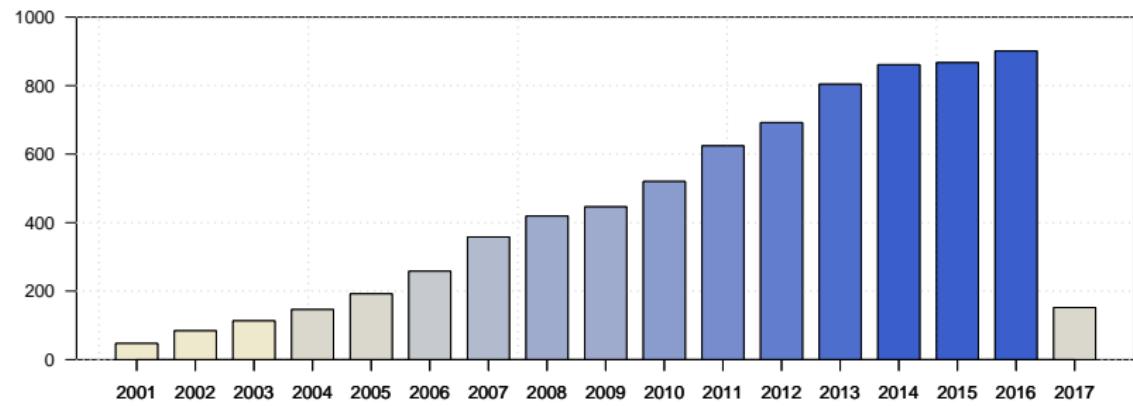


Figure 9: Citations per year of [Lee and Seung 1999](#)'s paper that launched NMF

References

-  Donoho, David L and Stodden, Victoria C (2004). 'When does non-negative matrix factorization give a correct decomposition into parts?' In: *Advances in neural information processing systems 16: proceedings of the 2003 conference*. Ed. by S Thrun, L Saul and B Schölkopf.
-  Hoyer, Patrik O (2004). 'Non-negative matrix factorization with sparseness constraints'. In: *Journal of machine learning research* 5.Nov, pp. 1457–1469.
-  Kim, Hyunsoo and Park, Haesun (2007). 'Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis'. In: *Bioinformatics* 23.12, pp. 1495–1502.
-  Lee, Daniel D and Seung, H Sebastian (1999). 'Learning the parts of objects by non-negative matrix factorization'. In: *Nature* 401.6755, pp. 788–791.
-  – (2001). 'Algorithms for non-negative matrix factorization'. In: *Advances in neural information processing systems*, pp. 556–562.
-  Paatero, Pentti and Tapper, Unto (1994). 'Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values'. In: *Environmetrics* 5.2, pp. 111–126.