# Topic modeling

Singular value decomposition (SVD)

& Non-negative matrix factorization (NMF)

# Singular value decomposition

- **Factorization of a real or complex matrix.**

- **Generalizes the eigendecomposition of a square normal matrix with an orthonormal eigenbasis to any $\Lambda$ matrix.**

- **SVD of an m × n matrix A is a factorization of the form**

$$A = U\Lambda V^T$$
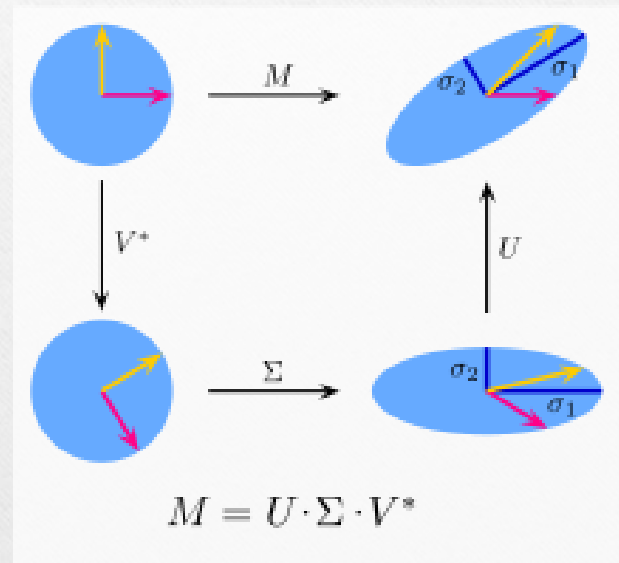
- U is an m × n unitary matrix, $UU^T = I$

- $\Lambda$ is an m × n rectangular diagonal matrix with non-negative real numbers on the diagonal,

- V is an n × n unitary matrix,

- $V^T$ is the conjugate transpose of V. $VV^T = I$

# Results

- $U$ is the eigenvector of $\mathbf{AA^T}$. Thus, $\mathbf{AA^T} = U(\Sigma\Sigma^T)U^T$ (left singular vector of $A$)

- $V$ is the eigenvector of $\mathbf{A^TA}$.

- Thus, $\mathbf{A^T A} = V(\Sigma^T\Sigma)V^T$ (right singular vector of $A$)

- $\sigma_i$ is the square root of the eigenvalues of $\mathbf{A^TA}$ or $\mathbf{AA^T}$.

$$A = \sum_{i=1}^{s} \sigma_i \boldsymbol{u_i v_i^T}$$

# Geometric interpretation of SVD


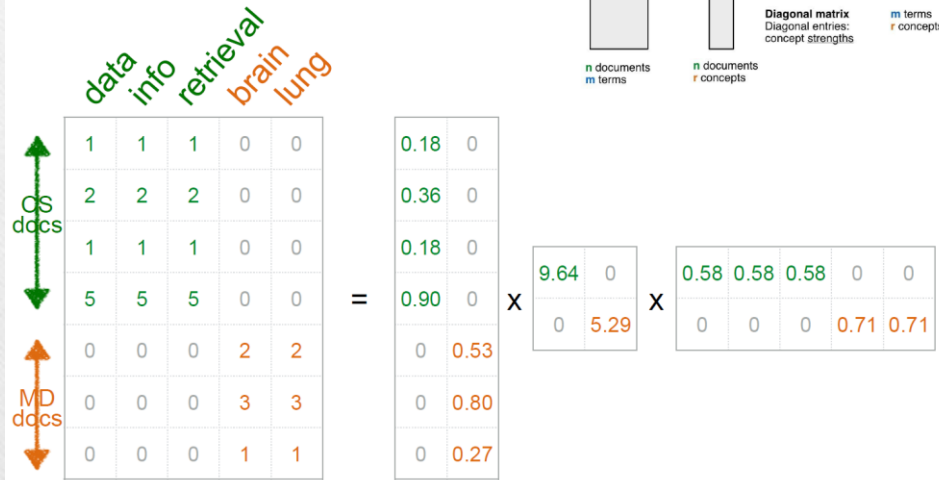
$$M = U \cdot \Sigma \cdot V^*$$

Source: https://en.wikipedia.org/wiki/Singular_value_decomposition

# R code

```
library(OpenImageR)

x<-readImage("C:/Users/user/Documents/pansy.jpeg")

r<-rgb_2gray(x)

imageShow(r)

 r.svd<-svd(r)

plot(r.svd$d)

u<-r.svd$u

v<-r.svd$v

d<-diag(r.svd$d)

 depth<-50

us<-as.matrix(u[,1:depth])

vs<-as.matrix(v[,1:depth])

ds<-as.matrix(d[1:depth,1:depth])

ls<-us%*%ds%*%t(vs)

imageShow(ls)
```
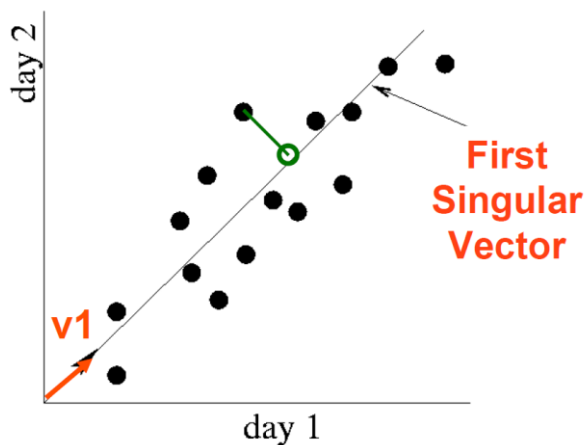
# Topic modeling using SVD

# SVD -Interpretation #1

- 'documents', 'terms' and 'concepts':

- **U**: document-concept similarity matrix
- **V**: term-concept similarity matrix
- $\Lambda$: diagonal elements: concept "strengths"

- $\mathbf{A^T}A =?$
- $\mathbf{AA^T} =?$

# SVD -Interpretation #2



**SVD is closely related to PCA**

$$
\begin{array}{|ccccc|}
\hline
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1 \\
\hline
\end{array}
=
\begin{array}{|cc|}
\hline
0.18 & 0 \\
0.36 & 0 \\
0.18 & 0 \\
0.90 & 0 \\
0 & 0.53 \\
0 & 0.80 \\
0 & 0.27 \\
\hline
\end{array}
\times
\begin{array}{|cc|}
\hline
9.64 & 0 \\
0 & 5.29 \\
\hline
\end{array}
\times
\begin{array}{|ccccc|}
\hline
0.58 & 0.58 & 0.58 & 0 & 0 \\
0 & 0 & 0 & 0.71 & 0.71 \\
\hline
\end{array}
$$

variance ('spread') on the v1 axis

v1

**A  =  U  $\Lambda$  $V^T$**

# How do we determine # of topics?

More details

Q: how exactly is dim. reduction done?

A: set the smallest singular values to zero:

# SVD -Interpretation #3

- finds non-zero 'blobs' in a data matrix =
- 'communities' (bi-partite cores, here)

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

Row 1
Row 4
Col 1
Col 3

Row 5
Row 7
Col 4

# NMF

- Given $X \in \mathbb{R}^{m \times n}$, compute an approximation $X \approx WH$ for some matrix $W \in \mathbb{R}^{m \times k}$ and $H \in \mathbb{R}^{k \times n}$ where $W$ and $H$ are nonnegative matrices

$$\min_{W \in \mathbb{R}^{m \times k}, \; H \in \mathbb{R}^{k \times n}} \|X - WH\|_F^2$$
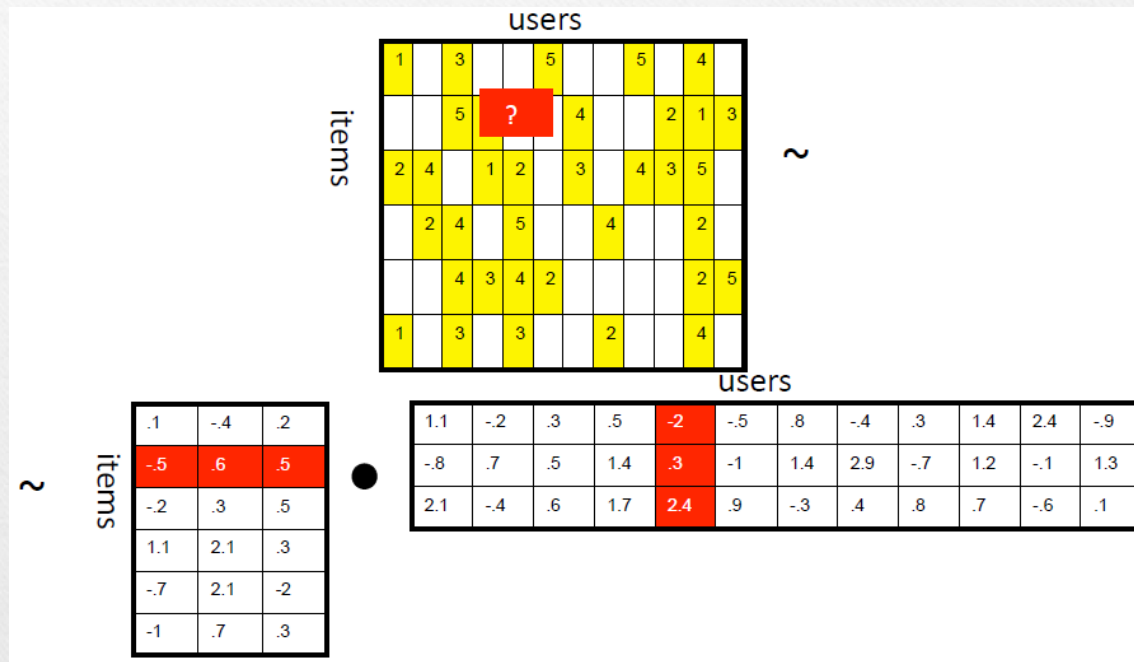
subject to

$$W \geq 0$$

$$H \geq 0$$

Often positive factors will be **more easily interpretable**

# Example

# Applications to Topic modeling