

Data Mining

Kapitel 10: Dimensionality Reduction

Johannes Zschache
Wintersemester 2018/19

Abteilung Datenbanken, Universität Leipzig
<http://dbs.uni-leipzig.de>

Übersicht

Hochdimension. Daten

Locality sensitive hashing

Clustering

Dimension. reduction

Graphdaten

PageRank, SimRank

Network Analysis

Spam Detection

Unbegrenzte Daten

Filtering data streams

Web advertising

Queries on streams

Maschinelles Lernen

Support Vector Machines

Decision Trees

Nearest Neighbors

Anwendung

Recommen. Systems

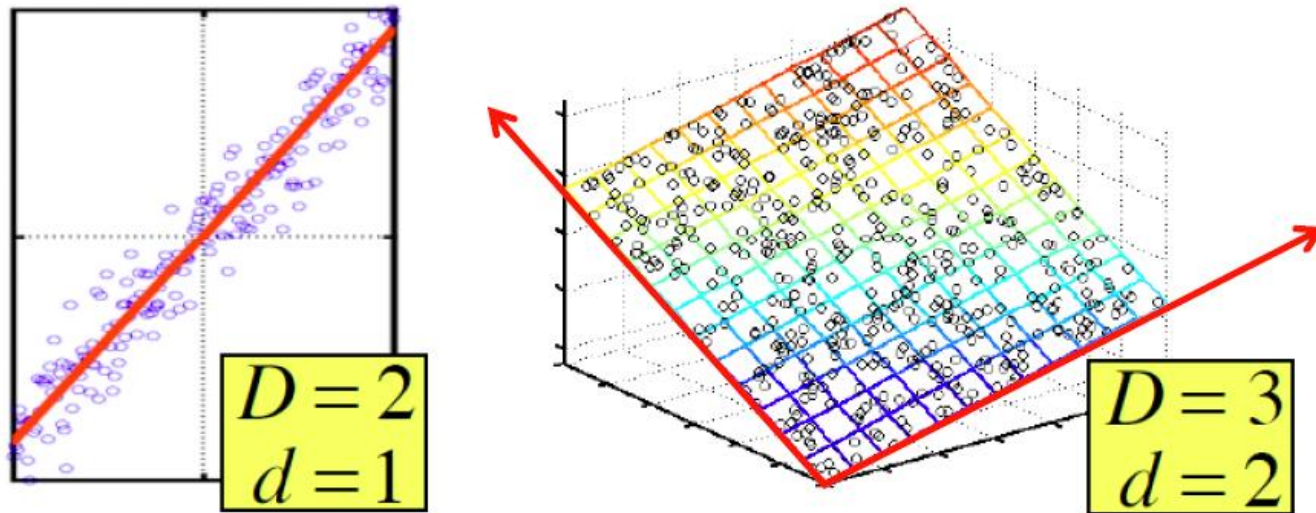
Association Rules

Duplicate document detection

Inhaltsverzeichnis

- **Einführung**
- **Singular Value Decomposition**
- **CUR Decomposition**

Dimensionsreduktion



- **Idee:** Falls Daten eines D -dimensionalen Raumes in der Nähe eines d -dimensionalen Unterraumes liegen, können die Achsen des Unterraums als Faktoren dienen
- Wahl der Faktoren:
 - Erster Faktor zeigt in die Richtung, in welcher die Punkte ihre größte Streuung aufweisen
 - Zweiter Faktor ist orthogonal zum ersten Faktor und zeigt in die Richtung mit der größten Streuung unter den Punkten, usw..

Dimensionsreduktion

- Einfaches Beispiel:

Kunde	Montag	Dienstag	Mittwoch	Donnerstag	Freitag	Repräsentation
A	1	1	1	0	0	[1, 0]
B	2	2	2	0	0	[2, 0]
C	1	1	1	0	0	[1, 0]
D	5	5	5	0	0	[5, 0]
E	0	0	0	2	2	[0, 2]
F	0	0	0	3	3	[0, 3]

- Rang** einer Matrix: Anzahl der *linear unabhängigen* Zeilen/Spalten

– Im Beispiel: Zeilenrang ist 2

– Aufspaltung:

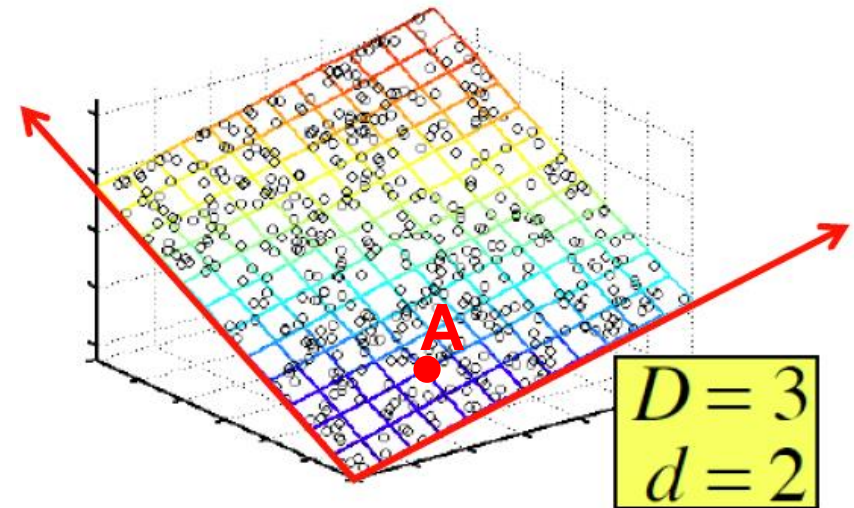
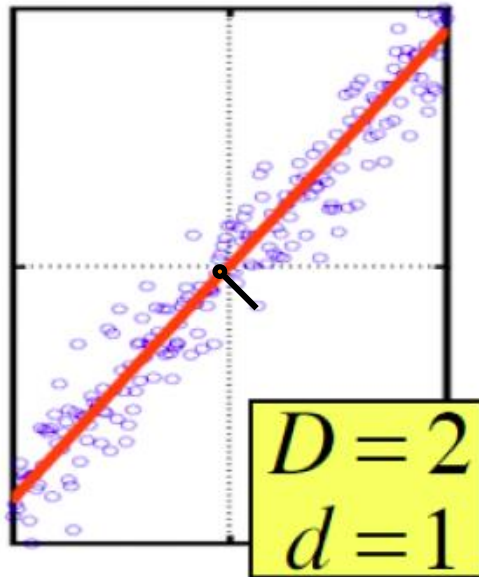
$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 0 \\ 1 & 0 \\ 5 & 0 \\ 0 & 2 \\ 0 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

Neue Achsen/Faktoren

Neue Koordinaten

Dimensionsreduktion

Ziel: Aufdeckung der Datenachsen



Anstatt 2 Koordinaten, wird jeder Punkt nur über eine Koordinate repräsentieren: Position auf der roten Linie

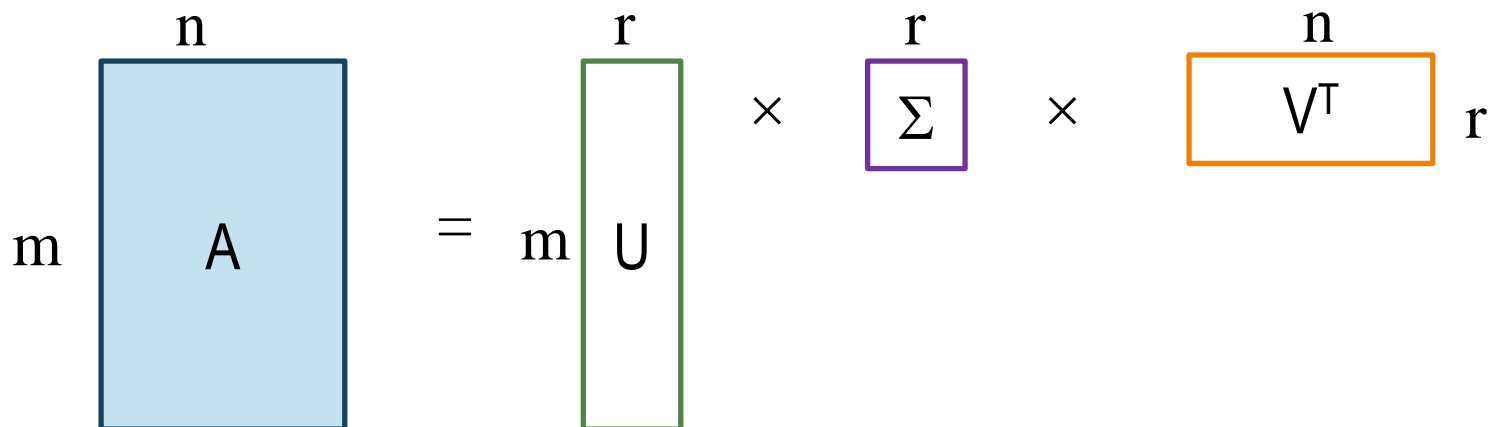
Punkt A wird z.B. anstatt durch $[3,4,2]$ durch $[2,1]$ repräsentiert (neue x-Achse ist Vektor $(1,2,1)$)

Inhaltsverzeichnis

- **Einführung**
- **Singular Value Decomposition**
- **CUR Decomposition**

Singular Value Decomposition (SVD)

- Zerlegung eine Matrix A in das Produkt dreier Matrizen:



$$A_{[m \times n]} = U_{[m \times r]} \cdot \Sigma_{[r \times r]} \cdot (V_{[n \times r]})^T$$

- Matrix Σ ist Diagonalmatrix deren Einträge ($\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$) nicht-negativ sind und als **Singulärwerte** bezeichnet werden
- Spalten von U und V sind orthonormal: $U^T \cdot U = V^T \cdot V = I$ (Einheitsmatrix)
- Anzahl der Faktoren = r = Rang von A

Beispiel: Nutzer und Filme

U ist Nutzer-Faktor-Matrix

Romantik \rightarrow \leftarrow SciFi \rightarrow

Nutzer

	Matrix	Alien	Serenity	Casablanca	Amelie
	1	1	1	0	0
	3	3	3	0	0
	4	4	4	0	0
	5	5	5	0	0
	0	2	0	4	4
	0	0	0	5	5
	0	1	0	2	2

SciFi-Faktor

Romantik-Faktor

Stärke des SciFi-Faktors

$$\begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} = \mathbf{X} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \mathbf{X}$$

V ist Film-Faktor-Matrix

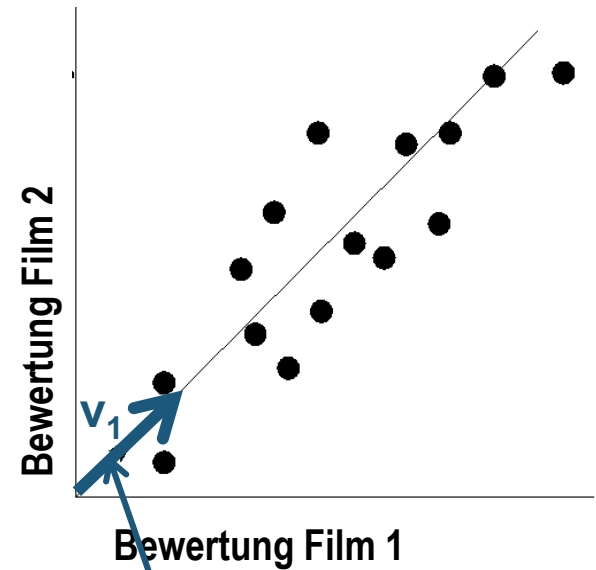
$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

SVD: Dimensionsreduktion

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \mathbf{X}$$

$$\begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \mathbf{X}$$

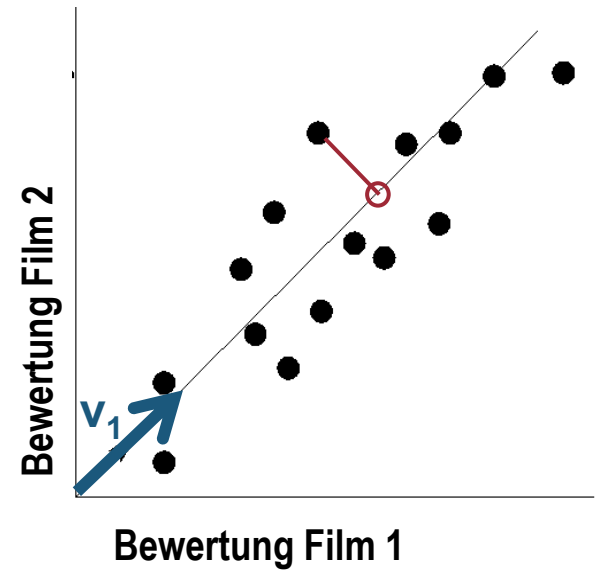
$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$



SVD: Dimensionsreduktion

$$\begin{bmatrix}
 \mathbf{0.13} & \mathbf{0.02} & \mathbf{-0.01} \\
 \mathbf{0.41} & \mathbf{0.07} & \mathbf{-0.03} \\
 \mathbf{0.55} & \mathbf{0.09} & \mathbf{-0.04} \\
 \mathbf{0.68} & \mathbf{0.11} & \mathbf{-0.05} \\
 \mathbf{0.15} & \mathbf{-0.59} & \mathbf{0.65} \\
 \mathbf{0.07} & \mathbf{-0.73} & \mathbf{-0.67} \\
 \mathbf{0.07} & \mathbf{-0.29} & \mathbf{0.32}
 \end{bmatrix}
 \mathbf{X}
 \begin{bmatrix}
 \mathbf{12.4} & \mathbf{0} & \mathbf{0} \\
 \mathbf{0} & \mathbf{9.5} & \mathbf{0} \\
 \mathbf{0} & \mathbf{0} & \mathbf{1.3}
 \end{bmatrix}
 \mathbf{T}
 =
 \begin{bmatrix}
 \mathbf{1.61} & \mathbf{0.19} & \mathbf{-0.01} \\
 \mathbf{5.08} & \mathbf{0.66} & \mathbf{-0.03} \\
 \mathbf{6.82} & \mathbf{0.85} & \mathbf{-0.05} \\
 \mathbf{8.43} & \mathbf{1.04} & \mathbf{-0.06} \\
 \mathbf{1.86} & \mathbf{-5.60} & \mathbf{0.84} \\
 \mathbf{0.86} & \mathbf{-6.93} & \mathbf{-0.87} \\
 \mathbf{0.86} & \mathbf{-2.75} & \mathbf{0.41}
 \end{bmatrix}$$

Projektionen auf die Achse
des SciFi-Faktors $(U \Sigma)^T$



SVD: Dimensionsreduktion

Reduktion der Dimensionen: Setze kleinsten Singulärwerte auf Null

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \mathbf{X} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{X} \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

Rang-2-Approximation von A (je größer der Rang desto genauer die Approximation)

SVD: Dimensionsreduktion

$$\underbrace{\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix}}_A \approx \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \mathbf{X} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{X}$$

$$\underbrace{\begin{bmatrix} 0.92 & 0.95 & 0.92 & 0.01 & 0.01 \\ 2.91 & 3.01 & 2.91 & -0.01 & -0.01 \\ 3.90 & 4.04 & 3.90 & 0.01 & 0.01 \\ 4.82 & 5.00 & 4.82 & 0.03 & 0.03 \\ 0.70 & 0.53 & 0.70 & 4.11 & 4.11 \\ -0.69 & 1.34 & -0.69 & 4.78 & 4.78 \\ 0.32 & 0.23 & 0.32 & 2.01 & 2.01 \end{bmatrix}}_B = \begin{bmatrix} 0.92 & 0.95 & 0.92 & 0.01 & 0.01 \\ 2.91 & 3.01 & 2.91 & -0.01 & -0.01 \\ 3.90 & 4.04 & 3.90 & 0.01 & 0.01 \\ 4.82 & 5.00 & 4.82 & 0.03 & 0.03 \\ 0.70 & 0.53 & 0.70 & 4.11 & 4.11 \\ -0.69 & 1.34 & -0.69 & 4.78 & 4.78 \\ 0.32 & 0.23 & 0.32 & 2.01 & 2.01 \end{bmatrix}$$

Genauigkeit über
Frobeniusnorm:

$$\|A - B\|_F$$

$$= \sqrt{\sum_{ij} (A_{ij} - B_{ij})^2}$$

SVD

- **Satz:** Sei k mit $0 \leq k \leq r$ die Anzahl der gewünschten Faktoren, $A = U \Sigma V^T$ und $B = U S V^T$ wobei S aus Σ konstruiert wurde, indem die letzten $r - k$ Diagonalelemente auf Null gesetzt wurden. Dann gilt:

$$B = \min_C \|A - C\|_F$$

- Für eine gegebene Anzahl an Faktoren k minimiert SVD den Fehler $\|A - B\|_F$, so dass B die *beste* Rang- k -Approximation für A darstellt
- Wie klein sollte man k wählen?
- Behalte 80-90% der „Energie“ $\sum_i \sigma_i^2$ (Summe über die quadrierten Diagonalelemente von Σ)
- Beispiel: Singulärwerte 12.4, 9.5, und 1.3 \rightarrow Energie: 245.7
 - Entfernen des letzten Singulärwertes setzt Energie auf 244 (99%)
 - Mit nur dem größten Singulärwert wäre die Energie auf 63% reduziert

SVD: Berechnung

- SVD für eine Matrix $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- Es gilt: $\mathbf{A}^T = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = (\mathbf{V}^T)^T \mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T$
 - Regel für die Transponierte eines Produkts von Matrizen
 - Zweifache Transposition löst sich auf
 - Transposition einer Diagonalmatrix ergibt die selbe Diagonalmatrix
- Somit gilt: $\mathbf{A}^T \mathbf{A} = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T$
 - Da Spalten von \mathbf{U} orthonormal: $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ (Identitätsmatrix)
 - $\mathbf{\Sigma}^2$ ist eine Diagonalmatrix deren i -tes Diagonalelement das Quadrat des i -ten Diagonalelements von $\mathbf{\Sigma}$ ist
- Da auch die Spalten von \mathbf{V} orthonormal: $\mathbf{A}^T \mathbf{A} \mathbf{V} = \mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T \mathbf{V} = \mathbf{V}\mathbf{\Sigma}^2$

D.h. \mathbf{V} ist die Matrix aus Eigenvektoren von $\mathbf{A}^T \mathbf{A}$ und die Diagonalelemente von $\mathbf{\Sigma}^2$ sind die dazugehörigen Eigenwerte.

- Analog: $\mathbf{A} \mathbf{A}^T \mathbf{U} = \mathbf{U}\mathbf{\Sigma}^2$

Eigenwerte und -vektoren

- Sei M eine quadratische Matrix. Eine reelle Zahl λ heißt Eigenwert von M und ein Vektor $e \neq 0$ der dazugehörige Eigenvektor, falls $Me = \lambda e$.
- Beispiel:

$$\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \cdot \begin{bmatrix} \frac{1}{\sqrt{5}} \\ 2 \\ \frac{1}{\sqrt{5}} \end{bmatrix} = 7 \begin{bmatrix} \frac{1}{\sqrt{5}} \\ 2 \\ \frac{1}{\sqrt{5}} \end{bmatrix}$$

- Somit ist $\begin{bmatrix} \frac{1}{\sqrt{5}} \\ 2 \\ \frac{1}{\sqrt{5}} \end{bmatrix} \approx \begin{bmatrix} 0.447 \\ 0.894 \end{bmatrix}$ ein Eigenvektor der Matrix $\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}$ und 7 der dazugehörige Eigenwert

- Berechnung der Eigenvektoren und -werte z.B. über **die Power-Iteration-Methode** (siehe auch Folie 4-14 aus „Link Analysis“)

Eigenwerte und -vektoren

- Power-Iteration-Methode:
 - Zu Beginn: beliebiger Vektor $x_0 \neq 0$
 - Iteration: $x_{k+1} = \frac{Mx_k}{\|Mx_k\|}$
 - $\| \dots \|$ bezeichnet die Frobeniusnorm
 - Stopp, falls Änderungen in x_k vernachlässigbar klein
- Beispiel: $M = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}$ und $x_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
 - $Mx_0 = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 8 \end{bmatrix}$ und $\|Mx_0\| = \sqrt{5^2 + 8^2} = 9.434$
 - $x_1 = \begin{bmatrix} 0.530 \\ 0.848 \end{bmatrix}$
 - $Mx_1 = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 0.530 \\ 0.848 \end{bmatrix} = \begin{bmatrix} 3.286 \\ 6.148 \end{bmatrix}$ und $\|Mx_1\| = 6.971$
 - $x_2 = \begin{bmatrix} 0.471 \\ 0.882 \end{bmatrix}$
 - ...

Eigenwerte und -vektoren

- Die Power-Iteration-Methode berechnet den ersten Eigenvektor x (mit dem größten Eigenwert)

- Dazugehörige Eigenwert: $\lambda = x^T M x$

- Beispiel:

$$\begin{bmatrix} 0.447 \\ 0.894 \end{bmatrix}^T \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 0.447 \\ 0.894 \end{bmatrix} = 6.993$$

- Reduzierung der Matrix M um den Anteil, der durch den ersten Eigenwert und –vektor generiert wird:

$$M^* := M - \lambda x x^T$$

- Power-Iteration-Method auf M^* berechnet den zweiten Eigenvektor von M (mit dem zweitgrößten Eigenwert von M)

- Beispiel:

$$\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} - 6.993 \begin{bmatrix} 0.447 \\ 0.894 \end{bmatrix} \begin{bmatrix} 0.447 \\ 0.894 \end{bmatrix}^T = \begin{bmatrix} 2.601 & -0.797 \\ -0.797 & 0.413 \end{bmatrix}$$

Beispiel: Nutzer und Filme

U ist Nutzer-Faktor-Matrix

$$\begin{array}{c} \text{Romantik} \longrightarrow \leftarrow \text{SciFi} \longrightarrow \\ \text{Nutzer} \end{array} \begin{array}{c} \text{Matrix} \\ \text{Alien} \\ \text{Serenity} \\ \text{Casablanca} \\ \text{Amelie} \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{array}{c} \text{SciFi-Faktor} \\ \text{Romantik-Faktor} \end{array} \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \mathbf{X} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \mathbf{X}$$

V ist Film-Faktor-Matrix

$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

Beispiel: Nutzer und Filme

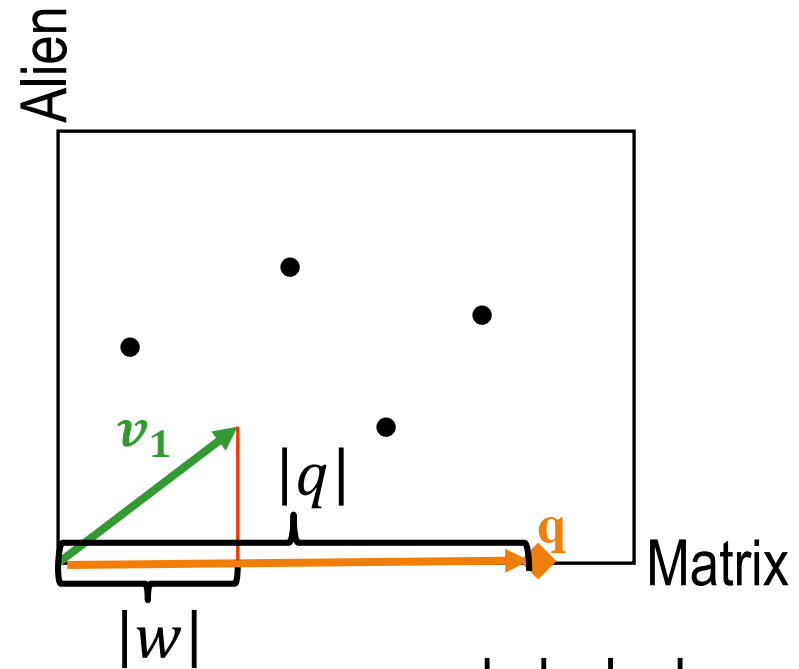
- Neuer Nutzer q mit Bewertung für Film "Matrix":

$$q = \begin{bmatrix} \text{Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ 5 & 0 & 0 & 0 & 0 \end{bmatrix}$$

- Projektion in den Film-Faktor-Raum

$$\begin{bmatrix} v_1 & v_2 \\ 0.56 & 0.12 \\ 0.59 & -0.02 \\ 0.56 & 0.12 \\ 0.09 & -0.69 \\ 0.09 & -0.69 \end{bmatrix}$$

Ersten beiden Spalten der
Film-Faktor-Matrix V



SciFi-Faktor $q \cdot v_1 = 2.8$

Romantik-Faktor $q \cdot v_2 = 0.6$

Beispiel: Nutzer und Filme

- Neuer Nutzer q mit Bewertung für Film "Matrix":

$$q = \begin{bmatrix} \text{Matrix} \\ 5 \\ \text{Alien} \\ 0 \\ \text{Serenity} \\ 0 \\ \text{Casablanca} \\ 0 \\ \text{Amelie} \\ 0 \end{bmatrix}$$

- Faktor-Repräsentation: $\begin{bmatrix} q \cdot v_1 & q \cdot v_2 \\ 2.8 & 0.6 \end{bmatrix}$

- Rückführung in den Film-Raum:

$$\begin{bmatrix} 2.8 & 0.6 \end{bmatrix} \mathbf{X} \begin{bmatrix} \mathbf{0.56} & \mathbf{0.59} & \mathbf{0.56} & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & \mathbf{-0.69} & \mathbf{-0.69} \end{bmatrix} =$$

$$\begin{bmatrix} \mathbf{1.64} & \mathbf{1.64} & \mathbf{1.64} & -1.62 & -1.62 \end{bmatrix}$$

Empfehlungen für Filme „Alien“ und „Serenity“

Beispiel: Nutzer und Filme

- **Beobachtung:** Ein weiterer Nutzer \mathbf{d} , der nur die Filme „Alien“ und „Serenity“ bewertet hat, kann, im Vergleich zum Nutzer \mathbf{q} , einen ähnlichen Faktor-Vektor haben; *obwohl es keinen Film gibt, der von beiden Nutzern bewertet wurde*

$$\begin{array}{rcc} & \begin{array}{c} \text{Matrix} \\ \text{Alien} \\ \text{Serenity} \\ \text{Casablanca} \\ \text{Amelie} \end{array} & \\ \mathbf{d} & = & \begin{bmatrix} 0 & 4 & 5 & 0 & 0 \end{bmatrix} \\ \mathbf{q} & = & \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \end{bmatrix} \end{array} \quad \begin{array}{c} \text{-----} \\ \text{-----} \end{array} \quad \begin{array}{cc} d \cdot v_1 & d \cdot v_2 \\ \left[\begin{array}{cc} 5.2 & 0.4 \end{array} \right] \\ q \cdot v_1 & q \cdot v_2 \\ \left[\begin{array}{cc} 2.8 & 0.6 \end{array} \right] \end{array}$$

Inhaltsverzeichnis

- **Einführung**
- **Singular Value Decomposition**
- **CUR Decomposition**

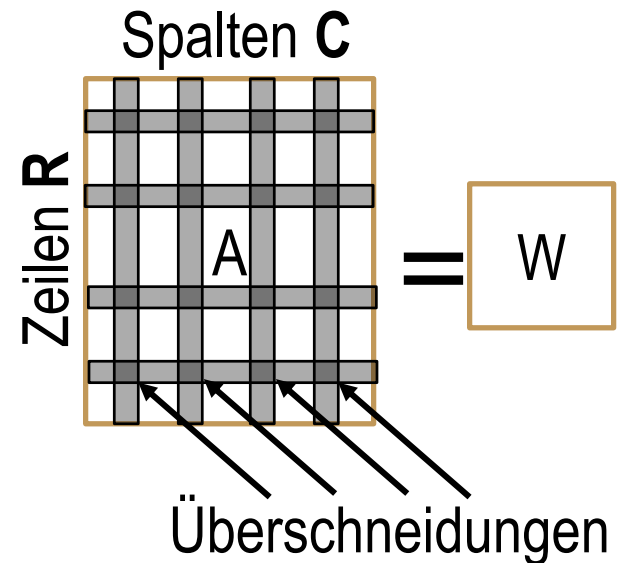
Berechnung der Matrix U

- W sei die Matrix aller Überschneidungen der Spalten C und der Reihen R
- Berechnung der SVD von W : $W = X Z Y^T$
- Dann ist $U := Y (Z^+)^2 X^T$

- Z^+ : Diagonalmatrix mit Diagonalelementen:

$$Z_{ii}^+ = \frac{1}{Z_{ii}} \text{ falls } Z_{ii} \neq 0 \text{ und } 0 \text{ sonst}$$

- Z^+ nennt man auch "Moore-Penrose-Pseudoinverse": Anstatt $ZZ^{-1} = I$ gilt $ZZ^+Z = Z$ und $Z^+ZZ^+ = Z^+$



- Beispiel:

$$W = \begin{bmatrix} 0 & 5 \\ 5 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Dann ist:

$$U = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1/5 & 0 \\ 0 & 1/5 \end{bmatrix}^2 \cdot \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1/25 \\ 1/25 & 0 \end{bmatrix}$$

Auswahl der Spalten und Zeilen

- Um den Fehler $\|A - C \cdot U \cdot R\|_F$ deutlich zu verringern sollten die Zeilen und Spalten nach *Wichtigkeit* ausgewählt werden
- Die Wichtigkeit einer Zeile/Spalte: quadrierte Frobeniusnorm
- Wahrscheinlichkeiten der Auswahl sind proportional zu deren Wichtigkeit
- Beispiel: Spalte [3,4,5] hat Wichtigkeit 50 und die Spalte [3,0,1] hat Wichtigkeit 10 → Wahrscheinlichkeit für erste Zeile ist fünfmal so groß wie Wahrscheinlichkeit der zweiten Zeile
- Algorithmus für Spalten: **Input:** matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, sample size c
Output: $\mathbf{C}_d \in \mathbb{R}^{m \times c}$
 1. for $x = 1 : n$ [column distribution]
 2. $P(x) = \sum_i \mathbf{A}(i, x)^2 / \sum_{i,j} \mathbf{A}(i, j)^2$
 3. for $i = 1 : c$ [sample columns]
 4. Pick $j \in 1 : n$ based on distribution $P(x)$
 5. Compute $\mathbf{C}_d(:, i) = \mathbf{A}(:, j) / \sqrt{cP(j)}$
- Eine Spalte kann mehrmals ausgewählt werden
- Skalierung der Spalten, um dies zu korrigieren

Beispiel

Wkt:									
0.012	1	1	1	0	0				
0.111	3	3	3	0	0				
0.198	4	4	4	0	0				
0.309	5	5	5	0	0				
0.132	0	0	0	4	4				
0.206	0	0	0	5	5				
0.033	0	0	0	2	2				

≈

<table style="border-collapse: collapse; text-align: center;"> <tr><td style="padding: 5px;">1.54</td><td style="padding: 5px;">0</td></tr> <tr><td style="padding: 5px;">4.63</td><td style="padding: 5px;">0</td></tr> <tr><td style="padding: 5px;">6.17</td><td style="padding: 5px;">0</td></tr> <tr><td style="padding: 5px;">7.72</td><td style="padding: 5px;">0</td></tr> <tr><td style="padding: 5px;">0</td><td style="padding: 5px;">6.58</td></tr> <tr><td style="padding: 5px;">0</td><td style="padding: 5px;">8.22</td></tr> <tr><td style="padding: 5px;">0</td><td style="padding: 5px;">3.29</td></tr> </table>	1.54	0	4.63	0	6.17	0	7.72	0	0	6.58	0	8.22	0	3.29	x	<table style="border-collapse: collapse; text-align: center;"> <tr> <td style="border-bottom: 1px solid black; padding: 5px;">0</td> <td style="padding: 5px;">1/25</td> </tr> <tr> <td style="padding: 5px;">1/25</td> <td style="padding: 5px;">0</td> </tr> </table>	0	1/25	1/25	0	x
1.54	0																				
4.63	0																				
6.17	0																				
7.72	0																				
0	6.58																				
0	8.22																				
0	3.29																				
0	1/25																				
1/25	0																				

=

<table style="border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 5px;">0.39</td> <td style="padding: 5px;">0.39</td> <td style="padding: 5px;">0.39</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> </tr> <tr> <td style="padding: 5px;">1.18</td> <td style="padding: 5px;">1.18</td> <td style="padding: 5px;">1.18</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> </tr> <tr> <td style="padding: 5px;">1.57</td> <td style="padding: 5px;">1.57</td> <td style="padding: 5px;">1.57</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> </tr> <tr> <td style="padding: 5px;">1.96</td> <td style="padding: 5px;">1.96</td> <td style="padding: 5px;">1.96</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> </tr> <tr> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">2.05</td> <td style="padding: 5px;">2.05</td> </tr> <tr> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">2.56</td> <td style="padding: 5px;">2.56</td> </tr> <tr> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">1.02</td> <td style="padding: 5px;">1.02</td> </tr> </table>	0.39	0.39	0.39	0	0	1.18	1.18	1.18	0	0	1.57	1.57	1.57	0	0	1.96	1.96	1.96	0	0	0	0	0	2.05	2.05	0	0	0	2.56	2.56	0	0	0	1.02	1.02	<table style="border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">7.79</td> <td style="padding: 5px;">7.79</td> </tr> <tr> <td style="padding: 5px;">6.36</td> <td style="padding: 5px;">6.36</td> <td style="padding: 5px;">6.36</td> <td style="padding: 5px;">0</td> <td style="padding: 5px;">0</td> </tr> </table>	0	0	0	7.79	7.79	6.36	6.36	6.36	0	0
0.39	0.39	0.39	0	0																																										
1.18	1.18	1.18	0	0																																										
1.57	1.57	1.57	0	0																																										
1.96	1.96	1.96	0	0																																										
0	0	0	2.05	2.05																																										
0	0	0	2.56	2.56																																										
0	0	0	1.02	1.02																																										
0	0	0	7.79	7.79																																										
6.36	6.36	6.36	0	0																																										

Wkt: 0.21

0.185

SVD vs. CUR

$$\text{SVD: } A = U \Sigma V^T$$

Klein und spärlich

Groß und spärlich

Groß und dicht

$$\text{CUR: } A = C U R$$

Klein und dicht

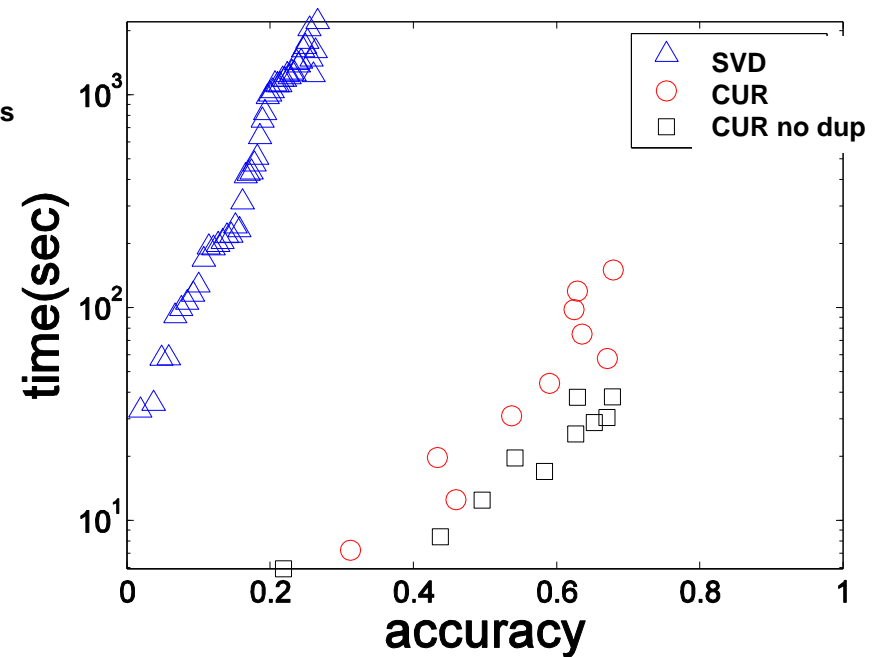
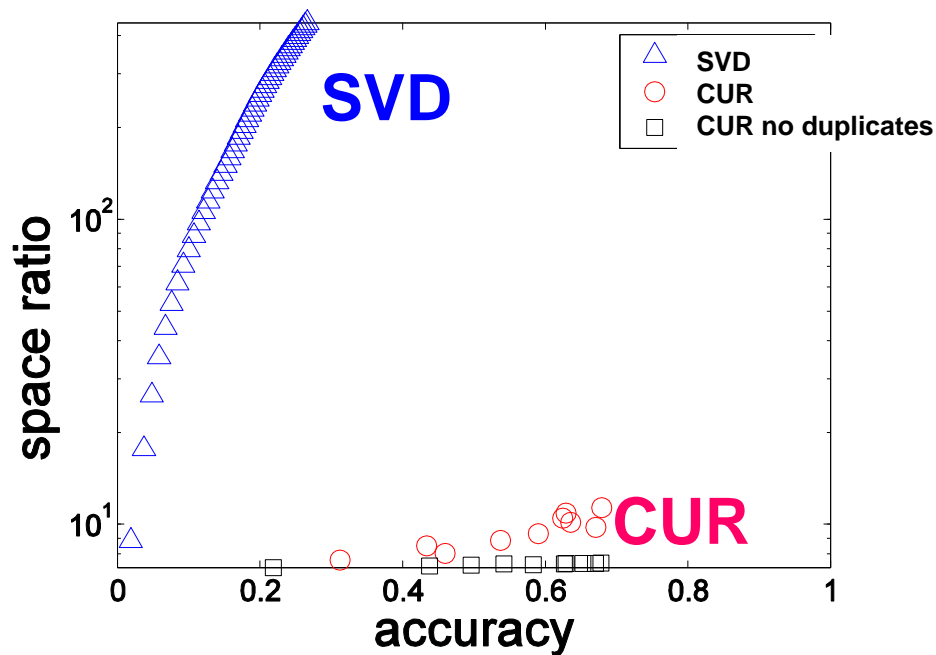
Groß und spärlich

Groß und spärlich

SVD vs. CUR: Experiment

- DBLP Daten
 - Bibliographische Sammlung wissenschaftlicher Publikationen im Bereich Informatik
 - Autor-Konferenz-Matrix A
 - Einträge A_{ij} : Anzahl der Publikationen des Autors i zur Konferenz j
 - 428 000 Autoren (Zeilen), 3 659 Konferenzen (Spalten)
 - Matrix ist sehr groß und spärlich besetzt
- Dimensionsreduktion über SVD und CUR
 - Wie lange laufen die Algorithmen?
 - Wie groß ist der Fehler zwischen approximierter und tatsächlicher Matrix
 - Wie viel Speicherplatz wird benötigt?

SVD vs. CUR: Experiment



- **Accuracy:** 1 – relative Summe der quadrierten Fehler
- **Space ratio:** Benötigter Speicherplatz
- **Time:** CPU Zeit

Sun, Faloutsos: *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM '07.

<http://www.cs.cmu.edu/~jimeng/papers/SunSDM07.pdf>

Referenzen, Beispiele, Übungen

Kapitel 11 aus „Mining of Massive Datasets“:

<http://www.mmids.org/>

Übungen:

- SVD: 11.3.1

Übung 11.3.1

$$M = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix}, M^T M = \begin{bmatrix} 36 & 37 & 38 \\ 37 & 49 & 61 \\ 38 & 61 & 84 \end{bmatrix}, M M^T = \begin{bmatrix} 14 & 26 & 22 & 16 & 22 \\ 26 & 50 & 46 & 28 & 40 \\ 22 & 46 & 50 & 20 & 32 \\ 16 & 28 & 20 & 20 & 26 \\ 22 & 40 & 32 & 26 & 35 \end{bmatrix}$$

Eigenvektoren $M^T M$:

$$\begin{bmatrix} -0.41 \\ -0.56 \\ -0.72 \end{bmatrix}, \begin{bmatrix} 0.82 \\ 0.13 \\ -0.56 \end{bmatrix}, \dots$$

Eigenvektoren $M M^T$:

$$\begin{bmatrix} -0.30 \\ -0.57 \\ -0.52 \\ -0.32 \\ -0.46 \end{bmatrix}, \begin{bmatrix} 0.16 \\ -0.03 \\ -0.74 \\ 0.51 \\ 0.41 \end{bmatrix}, \dots$$

Eigenwerte:

$$153.6, \quad 15.4, \quad 0.0, \dots$$

Übung 11.3.1

$$M = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix} = \underbrace{\begin{bmatrix} -0.30 & 0.16 \\ -0.57 & -0.03 \\ -0.52 & -0.74 \\ -0.32 & 0.51 \\ -0.46 & 0.41 \end{bmatrix}}_U \cdot \underbrace{\begin{bmatrix} 12.4 & 0 \\ 0 & 3.9 \end{bmatrix}}_{\Sigma} \cdot \underbrace{\begin{bmatrix} -0.41 & -0.56 & -0.72 \\ 0.82 & 0.13 & -0.56 \end{bmatrix}}_{V^T}$$

$$\begin{bmatrix} -0.30 & 0.16 \\ -0.57 & -0.03 \\ -0.52 & -0.74 \\ -0.32 & 0.51 \\ -0.46 & 0.41 \end{bmatrix} \cdot \begin{bmatrix} 12.4 & 0 \\ 0 & \mathbf{0} \end{bmatrix} \cdot \begin{bmatrix} -0.41 & -0.56 & -0.72 \\ 0.82 & 0.13 & -0.56 \end{bmatrix} = \begin{bmatrix} 1.5 & 2.1 & 2.6 \\ 2.9 & 4.0 & 5.1 \\ 2.6 & 3.6 & 4.6 \\ 1.6 & 2.3 & 2.9 \\ 2.3 & 3.2 & 4.1 \end{bmatrix}$$

Energie: $12.4^2 \approx 153.6$ und $12.4^2 + 3.9^2 = 169$, d.h. ca 91%