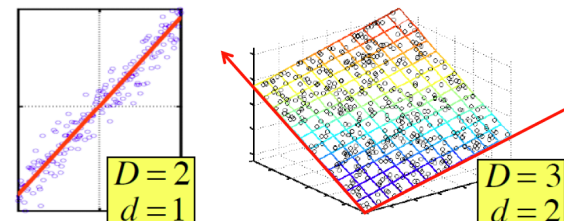


大数据分析

Scalable Machine Learning
Dimension Reduction

刘盛华

降维



- **Assumption:** 数据分布于低维度子空间
- 该子空间的基可以有效表征数据

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets,
<http://www.mmds.org>

2

矩阵的秩

- **Q:** 矩阵 A 的秩(rank)是什么?
- **A:** A 中线性无关的列的数量
- For example:
 - $A = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$ 秩 $r=2$
 - 原因: 1) A 的前两行线性无关, 所以秩至少为 2 2) A 第三行与前两行线性相关 ($\text{row}_3 = \text{row}_1 - \text{row}_2$), 所以秩小于 3
- **为什么需要关心矩阵的低秩?**
 - 矩阵 A 可以先写成两个“基”向量: $[1 \ 2 \ 1] \ [-2 \ -3 \ 1]$
 - 然后通过二维坐标表征: $[1 \ 0] \ [0 \ 1] \ [1 \ 1]$

3

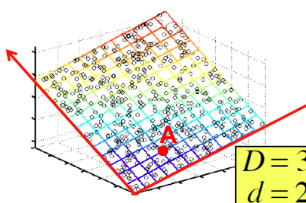
秩就是“维度”

- **3D 空间中的点云:**

- 可以用矩阵表征点的位置:

$$\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} \begin{matrix} A \\ B \\ C \end{matrix}$$

每行表示一个点:

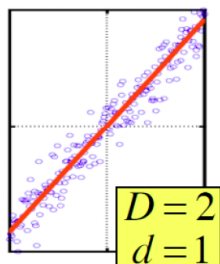


- **可以更简洁地重写每个点的坐标**
 - 旧的基向量: $[1 \ 0 \ 0] \ [0 \ 1 \ 0] \ [0 \ 0 \ 1]$
 - 新的基向量: $[1 \ 2 \ 1] \ [-2 \ -3 \ 1]$
 - 三个点都有新的坐标 A: $[1 \ 0]$, B: $[0 \ 1]$, C: $[1 \ -1]$
 - 注: 我们已经减少了坐标的维度

4

降维

- 降维的目的是寻找数据的轴(axis)



我们可以根据数据点在红线上的分布, 使用一维坐标替代原始的二维坐标。

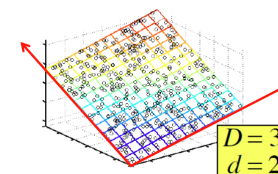
该方法会引入一定的偏差, 因为数据并没有严格分布在红线上。

5

为什么降维?

为什么降维?

- 发现隐含的关联、主题
 - 通常共现的词
- 移除冗余的、有噪声的特征
 - 并不是所有的词都有用
- 数据的解释、可视化
- 方便数据的存储和处理



6

SVD - Definition

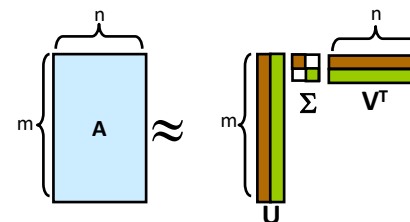
$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

- A: 输入矩阵
 - $m \times n$ matrix (e.g., m 篇文档, n 个词项 terms)
- U: A的左奇异向量构成的矩阵
 - $m \times r$ matrix (m 篇文档, r 个概念 concepts)
- Σ : 奇异值矩阵
 - $r \times r$ 对角阵 (每种主题的“强度” strength)
 - (r : 矩阵A的秩)
- V: A的右奇异向量构成的矩阵
 - $n \times r$ matrix (n 个词项 terms, r 个概念 concepts)

7

SVD

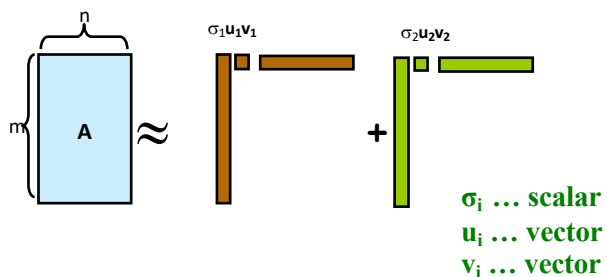
$$A \approx U \Sigma V^T = \sum_i \sigma_i u_i \circ v_i^T$$



8

SVD

$$A \approx U \Sigma V^T = \sum_i \sigma_i u_i \circ v_i^T$$



9

SVD - Properties

我们总能对一个实值矩阵A分解： $A = U \Sigma V^T$

分解时需要满足：

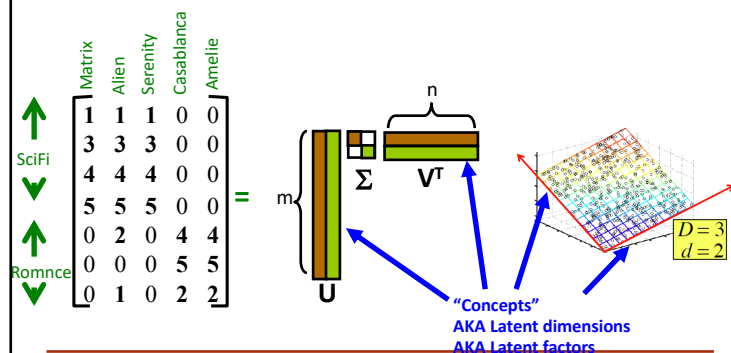
- U, Σ, V : **唯一**
- U, V : **每个矩阵的列向量彼此正交**
 - $U^T U = I; V^T V = I$ (I : 单位阵)
 - 每个矩阵的列都是彼此正交的单位向量
- Σ : **对角矩阵**
 - 每个元素 (奇异值) 是**正的**, 并且按照值大小降序排列 ($\sigma_1 \geq \sigma_2 \geq \dots \geq 0$)

Nice proof of uniqueness: <http://www.mpi-inf.mpg.de/~bast/ir-seminar-ws04/lecture2.pdf>

10

SVD – Example: Users-to-Movies

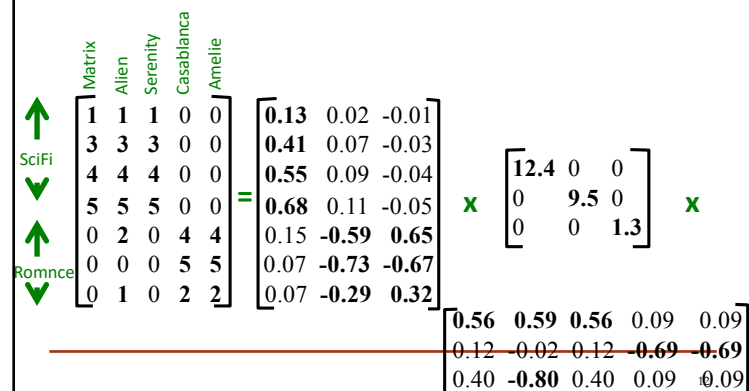
■ $A = U \Sigma V^T$ - example: 用户—电影评分矩阵



11

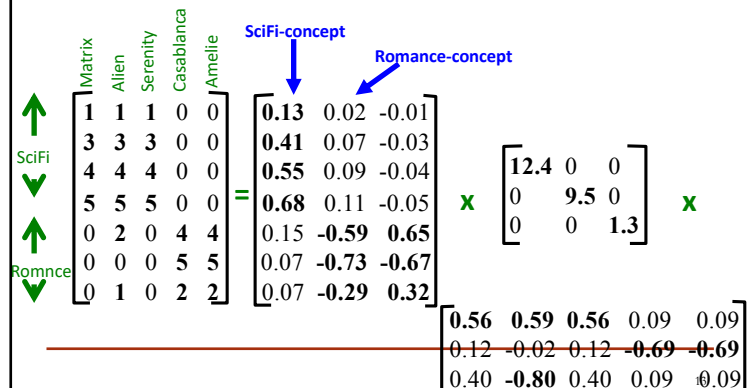
SVD – Example: Users-to-Movies

■ $A = U \Sigma V^T$ - example: 用户—电影评分矩阵



SVD – Example: Users-to-Movies

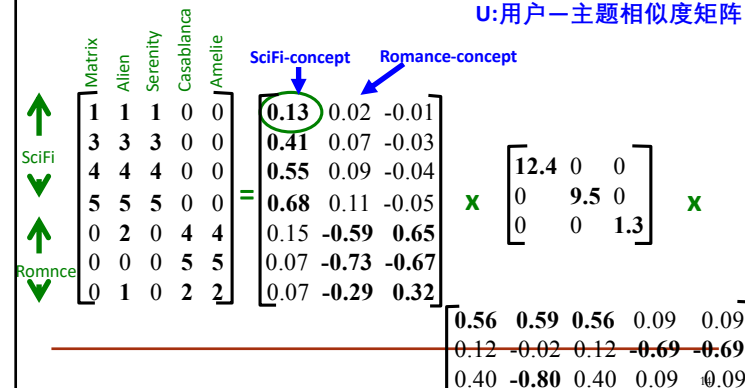
■ $A = U \Sigma V^T$ - example: 用户—电影评分矩阵



SVD – Example: Users-to-Movies

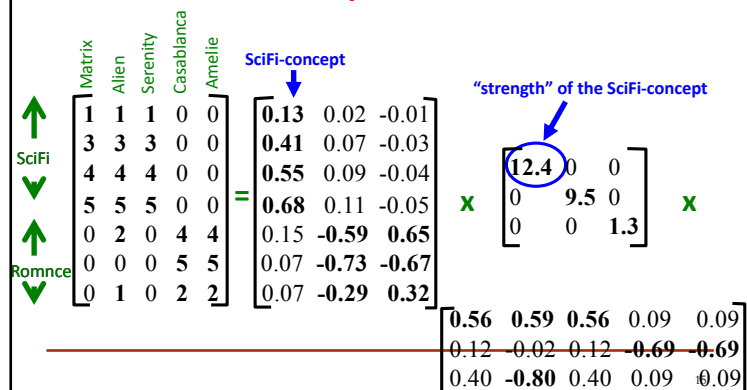
■ $A = U \Sigma V^T$ - example:

U: 用户—主题相似度矩阵



SVD – Example: Users-to-Movies

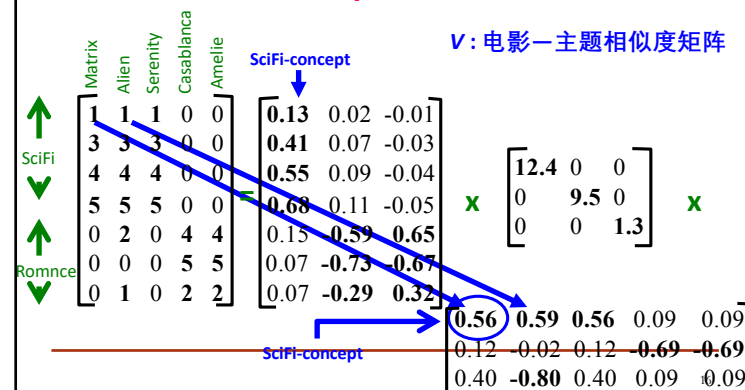
■ $A = U \Sigma V^T$ - example:



SVD – Example: Users-to-Movies

■ $A = U \Sigma V^T$ - example:

V: 电影—主题相似度矩阵



SVD – 解释 #1

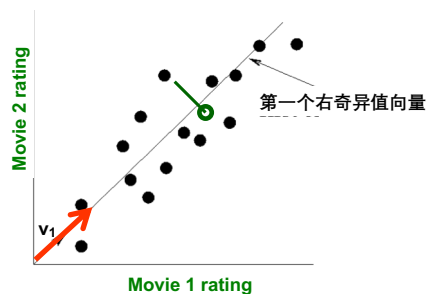
‘电影’, ‘用户’ and ‘主题’:
‘movies’, ‘users’ and ‘concepts’:

- U : 用户—主题相似度矩阵
- V : 电影—主题相似度矩阵
- Σ : 对角线的元素代表了每个主题**的强度**

17

使用SVD降维

SVD – 降维



- 只使用一维坐标 (z) 替代原始的二维坐标 (x, y)
- 点的坐标是其在向量 v_1 上的位置
- 如何选择 v_1 ? **最小化重构误差**

19

SVD降维

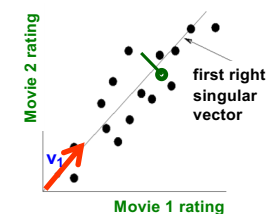
- Goal: **最小化总重构误差:**

$$\sum_{i=1}^N \sum_{j=1}^D \|x_{ij} - z_{ij}\|^2$$

- x_{ij} 是旧坐标 z_{ij} 是新坐标

- SVD给出最适合原始数据投影的轴:

- 最适合: 最小化重构误差
- 换言之, SVD给出最小的重构误差



20

SVD – 解释 #2

■ $A = U \Sigma V^T$ - example:

□ V : “电影—主题” 矩阵

□ U : “用户—主题” 矩阵

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & -0.09 \end{bmatrix}$$

SVD - Interpretation #2

■ $A = U \Sigma V^T$ - example:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & -0.09 \end{bmatrix}$$

variance ('spread') on the v_1 axis

SVD - Interpretation #2

$A = U \Sigma V^T$ - example:

■ $U \Sigma$: 给出投影轴以及数据点的投影值 (坐标)

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} = \begin{bmatrix} 1.61 & 0.19 & -0.01 \\ 5.08 & 0.66 & -0.03 \\ 6.82 & 0.85 & -0.05 \\ 8.43 & 1.04 & -0.06 \\ 1.86 & -5.60 & 0.84 \\ -0.86 & -6.93 & -0.87 \\ -0.86 & -2.75 & 0.41 \end{bmatrix}$$

Projection of users on the “Sci-Fi” axis ($U \Sigma$):

SVD - Interpretation #2

More details

■ Q: SVD降维有多精确?

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & -0.09 \end{bmatrix}$$

SVD – 最优低秩近似.

Theorem:

Let $A = U \Sigma V^T$ and $B = U S V^T$ where

S = diagonal $r \times r$ matrix with $s_i = \sigma_i$ ($i=1 \dots k$) else $s_i=0$
then B is a best rank(B)= k approx. to A

What do we mean by “best”:

B is a solution to $\min_B \|A-B\|_F$ where $\text{rank}(B)=k$

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} u_{11} & \dots & u_{1r} \\ \vdots & \ddots & \vdots \\ u_{m1} & \dots & u_{mr} \end{pmatrix} \begin{pmatrix} \sigma_{11} & 0 & \dots \\ 0 & \sigma_{22} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \vdots \\ v_{r1} & \dots & v_{rn} \end{pmatrix}$$

$$\|A-B\|_F = \sqrt{\sum_{ij} (A_{ij} - B_{ij})^2}$$

SVD - Interpretation #2

More details

- Q: SVD降维有多精确?
- A: 把最小的奇异值(一个或多个)设置为0.

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

SVD - Interpretation #2

More details

- Q: SVD降维有多精确?
- A: 把最小的奇异值(一个或多个)设置为0.

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

SVD - Interpretation #2

More details

- Q: SVD降维有多精确?
- A: 把最小的奇异值(一个或多个)设置为0.

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

SVD - Interpretation #2

More details

- Q: SVD降维有多精确?
- A: 把最小的奇异值(一个或多个)设置为0.

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & 0.02 \\ 0.41 & 0.07 \\ 0.55 & 0.09 \\ 0.68 & 0.11 \\ 0.15 & -0.59 \\ 0.07 & -0.73 \\ 0.07 & -0.29 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \end{bmatrix}$$

29

SVD - Interpretation #2

More details

- Q: SVD降维有多精确?
- A: 把最小的奇异值(一个或多个)设置为0.

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.92 & 0.95 & 0.92 & 0.01 & 0.01 \\ 2.91 & 3.01 & 2.91 & -0.01 & -0.01 \\ 3.90 & 4.04 & 3.90 & 0.01 & 0.01 \\ 4.82 & 5.00 & 4.82 & 0.03 & 0.03 \\ 0.70 & 0.53 & 0.70 & 4.11 & 4.11 \\ -0.69 & 1.34 & -0.69 & 4.78 & 4.78 \\ 0.32 & 0.23 & 0.32 & 2.01 & 2.01 \end{bmatrix}$$

Frobenius norm:

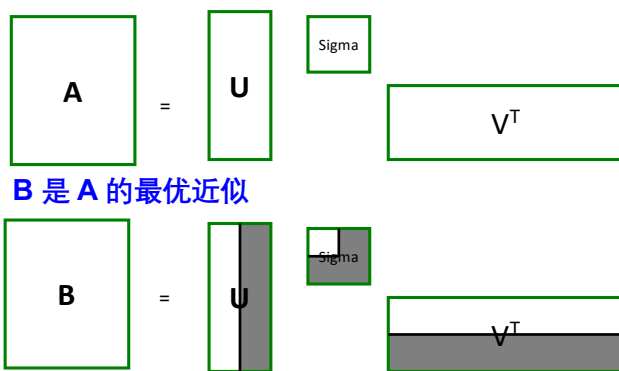
$$\|M\|_F = \sqrt{\sum_{ij} M_{ij}^2}$$

$$\|A-B\|_F = \sqrt{\sum_{ij} (A_{ij}-B_{ij})^2}$$

is "small"

30

SVD – 最优低秩近似.



31

SVD - Interpretation #2

等价于:
矩阵的谱分解

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} | & | \\ u_1 & u_2 \\ | & | \end{bmatrix} \times \begin{bmatrix} \sigma_1 & \emptyset \\ \emptyset & \sigma_2 \end{bmatrix} \times \begin{bmatrix} \text{---} v_1 & \text{---} \\ \text{---} v_2 & \text{---} \end{bmatrix}$$

32

SVD - Interpretation #2

等价于:
矩阵的谱分解

$$\begin{array}{c} \uparrow \\ n \\ \downarrow \end{array} \begin{array}{c} \leftarrow m \rightarrow \\ \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \end{array} = \begin{array}{c} \leftarrow k \text{ terms} \rightarrow \\ \sigma_1 \begin{array}{c} \nearrow u_1 \\ n \times 1 \end{array} \begin{array}{c} \nwarrow v_1^T \\ 1 \times m \end{array} + \sigma_2 \begin{array}{c} \nearrow u_2 \\ n \times 1 \end{array} \begin{array}{c} \nwarrow v_2^T \\ 1 \times m \end{array} + \dots \\ \text{Assume: } \sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0 \end{array}$$

为什么把最小的奇异值 σ_i 设置为0是正确的?
向量 u_i, v_i 是单位向量, 所以 σ_i 起到缩放的作用。把小的奇异值 σ_i 设为0只会引入少量误差。

33

SVD - Interpretation #2

Q: 保留多少奇异值 σ_i ?

A: 经验:

keep 80-90% of 'energy' = $\sum_i \sigma_i^2$

$$\begin{array}{c} \uparrow \\ n \\ \downarrow \end{array} \begin{array}{c} \leftarrow m \rightarrow \\ \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \end{array} = \begin{array}{c} \sigma_1 \begin{array}{c} \nearrow u_1 \\ n \times 1 \end{array} \begin{array}{c} \nwarrow v_1^T \\ 1 \times m \end{array} + \sigma_2 \begin{array}{c} \nearrow u_2 \\ n \times 1 \end{array} \begin{array}{c} \nwarrow v_2^T \\ 1 \times m \end{array} + \dots \\ \text{Assume: } \sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \end{array}$$

34

SVD - 复杂度分析

■ 计算SVD:

- $O(nm^2)$ or $O(n^2m)$ (whichever is less)

■ 但是:

- 可以通过更少的计算得到
 - 如果只需要奇异值,
 - 或者只需要前 k 个特征向量
 - 或者矩阵是稀疏的

■ 已经集成在许多线性代数计算包中:

- LINPACK, Matlab, Python, SPlus, Mathematica ...

35

SVD - 总结

■ SVD: $A = U \Sigma V^T$: 分解方法是唯一的

- U : 用户—主题相似度
- V : 电影—主题相似度
- Σ : 每个主题的程度

■ 降维:

- 保存几个最大的奇异值 (80-90% of 'energy')
- SVD: 通过数据的轴的线性组合重构数据

36