



中国科学院大学

University of Chinese Academy of Sciences

# Mining Massive Datasets

## Dimensionality Reduction

### SVD&CUR



# Dimensionality Reduction

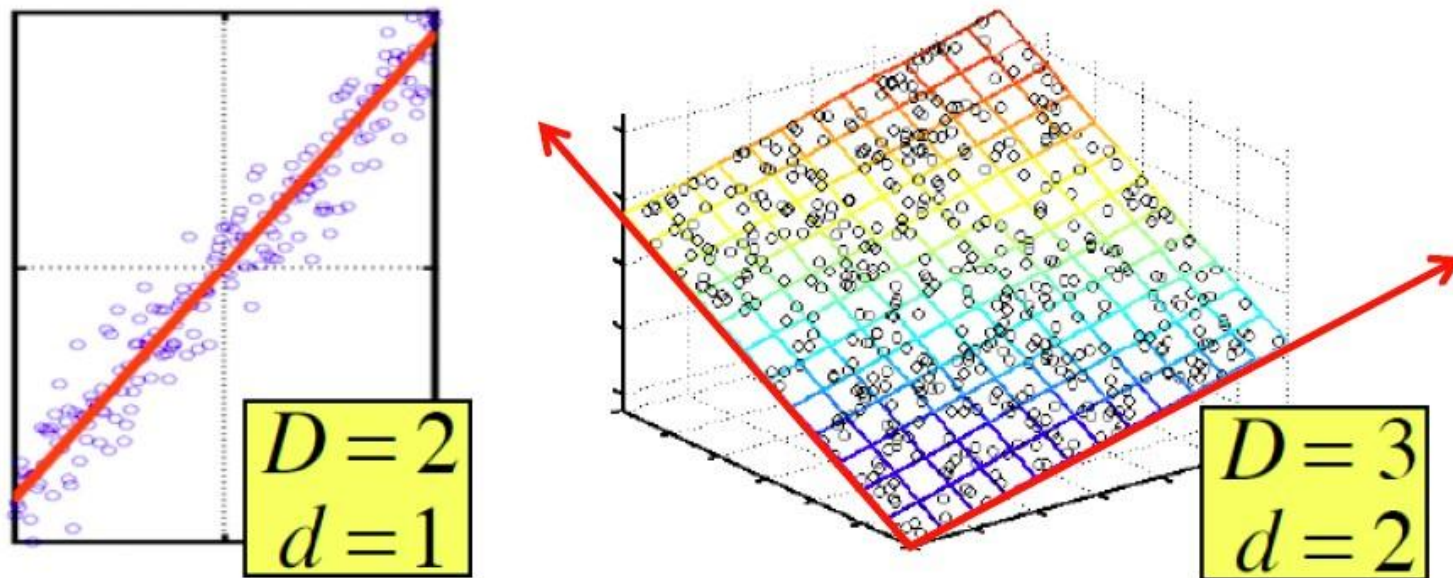
- **Compress / reduce dimensionality:**
  - $10^6$  rows;  $10^3$  columns; no updates
  - Random access to any cell(s); **small error: OK**

customer	day	We 7/10/96	Th 7/11/96	Fr 7/12/96	Sa 7/13/96	Su 7/14/96
ABC Inc.		1	1	1	0	0
DEF Ltd.		2	2	2	0	0
GHI Inc.		1	1	1	0	0
KLM Co.		5	5	5	0	0
Smith		0	0	0	2	2
Johnson		0	0	0	3	3
Thompson		0	0	0	1	1

The above matrix is really “2-dimensional.” All rows can be reconstructed by scaling  $[1 \ 1 \ 1 \ 0 \ 0]$  or  $[0 \ 0 \ 0 \ 1 \ 1]$



# Dimensionality Reduction



- **Assumption:** Data lies on or near a low  $d$ -dimensional subspace
- Axes of this subspace are effective representation of the data

# Why Reduce Dimensions?

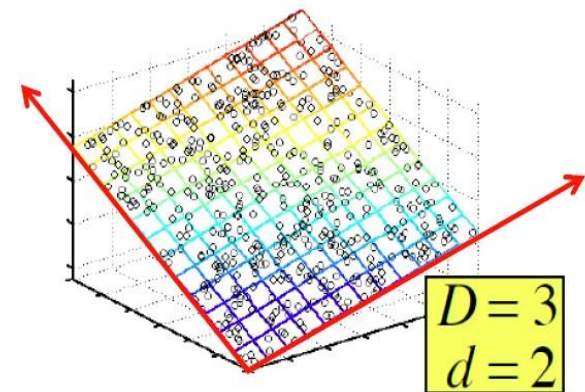
- Some features may be irrelevant
- We want to visualize high dimensional data
- “Intrinsic” dimensionality may be smaller than the number of features
- In particular, choose projection that minimizes the squared error in reconstructing original data



# Why Reduce Dimensions?

## Why reduce dimensions?

- Discover hidden correlations/topics
  - Words that occur commonly together
- Remove redundant and noisy features
  - Not all words are useful
- Interpretation and visualization
- Easier storage and processing of the data



# SVD-Definition

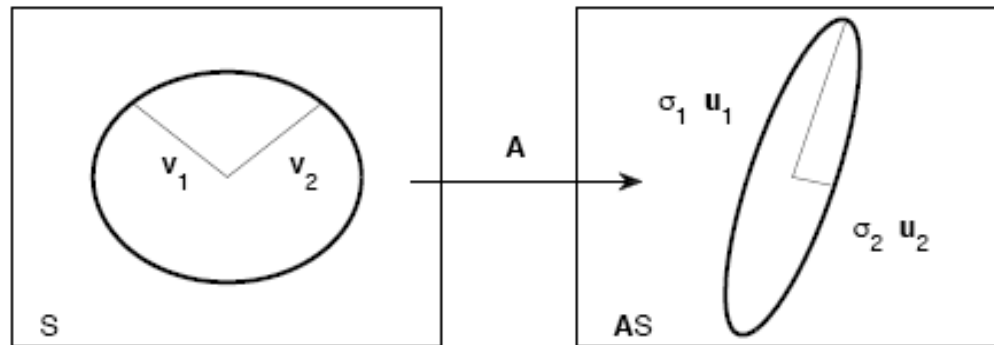
$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

- **A: Input data matrix**
  - $m \times n$  matrix (e.g.,  $m$  documents,  $n$  terms)
- **U: Left singular vectors**
  - $m \times r$  matrix ( $m$  documents,  $r$  concepts)
- **$\Sigma$ : Singular values**
  - $r \times r$  diagonal matrix (strength of each 'concept') ( $r$ : rank of the matrix **A**)
- **V: Right singular vectors**
  - $n \times r$  matrix ( $n$  terms,  $r$  concepts)



# SVD-decomposition

- The SVD, much as illustrated in the following figure, is essentially a transformation that stretches/compresses and rotates a given set of vectors



the transformation from the unit sphere to the hyperellipse

$$\mathbf{A}\mathbf{v}_j = \sigma_j \mathbf{u}_j,$$
$$\mathbf{A}^T \mathbf{u}_i = \sigma_j \mathbf{v}_j.$$



# SVD–Eigenvalue & Eigenvector

Given a  $n \times n$  matrix  $A^T A$ , for any  $\sigma$  and  $v$ , if

$$A^T A v_j = \sigma_j v_j$$

Then  $\sigma$  is called eigenvalue, and  $w$  is called eigenvector.

- To gain insight into the SVD, treat the rows of an  $m \times n$  matrix  $A$  as  $n$  points in a  $n$ -dimensional space and consider the problem of finding the best  $r$ -dimensional subspace with respect to the set of points. Here best means minimize the sum of the squares of the perpendicular distances of the points to the subspace.





# SVD-decomposition

- The objective of the rotation transformation is to find the maximal variance. We Projection of data along  $v$  is  $Av$ . Variance:

$$\sigma^2 = (Av)^T (Av) = v^T A^T Av$$

where  $A^T A$  is the covariance matrix of the data

Objective: maximize variance subject to constraint  $v^T v = 1$ .

Maximize  $f = v^T A^T Av - \lambda(v^T v - 1)$

$\lambda$  is the Lagrange multiplier, Differentiating with respect to  $v$  yields Eigenvalue equation:

$$A^T Av = \lambda v$$



# SVD-decomposition

- The objective of the rotation transformation is to find the maximal variance. We Projection of data along  $v$  is  $Av$ .

Variance:

$$\sigma^2 = (Av)^T (Av) = v^T A^T Av$$

where  $A^T A$  is the covariance matrix of the data

Objective: maximize variance subject to constraint  $v^T v = 1$ .

Maximize  $f = v^T A^T Av - \sigma(v^T v - 1)$

$\sigma$  is the Lagrange multiplier, Differentiating with respect to  $v$  yields Eigenvalue equation:

$$A^T Av = \lambda v$$



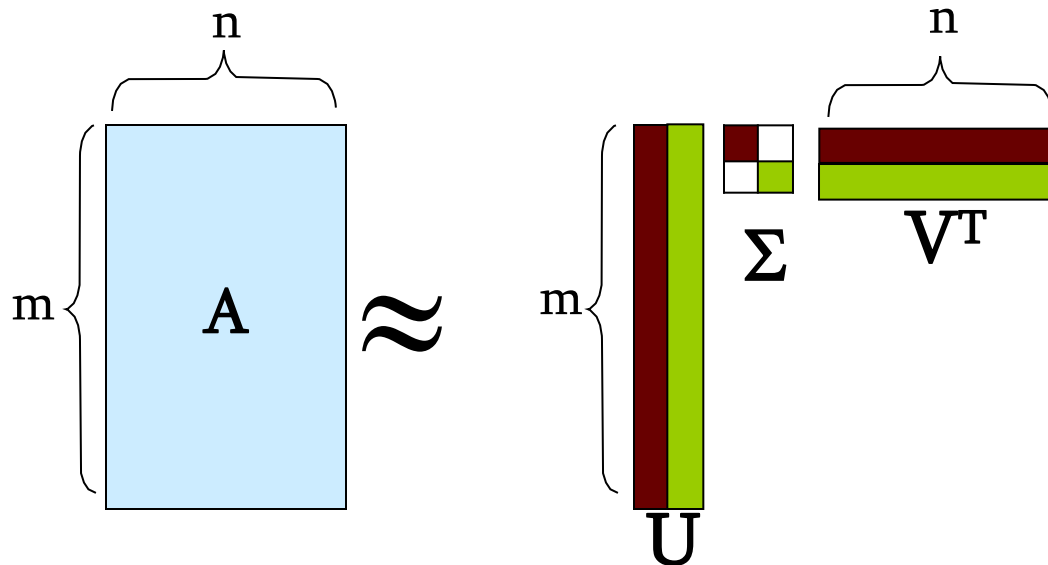
# SVD and PCA

- For symmetric  $A$ , SVD is closely related to PCA
- PCA:  $A = U \sigma U^T$ 
  - $U$  and  $\sigma$  are eigenvectors and eigenvalues.
- SVD:  $A = U \sigma V^T$ 
  - $U$  is left(column) eigenvectors
  - $V$  is right(row) eigenvectors
  - $\sigma$  is the same eigenvalues
- Note the difference of  $A$  in PCA and SVD
  - SVD:  $A$  is directly the data, e.g. word-by-document matrix
  - PCA:  $A$  is covariance matrix,  $A = X^T X$ , each row in  $X$  is a sample



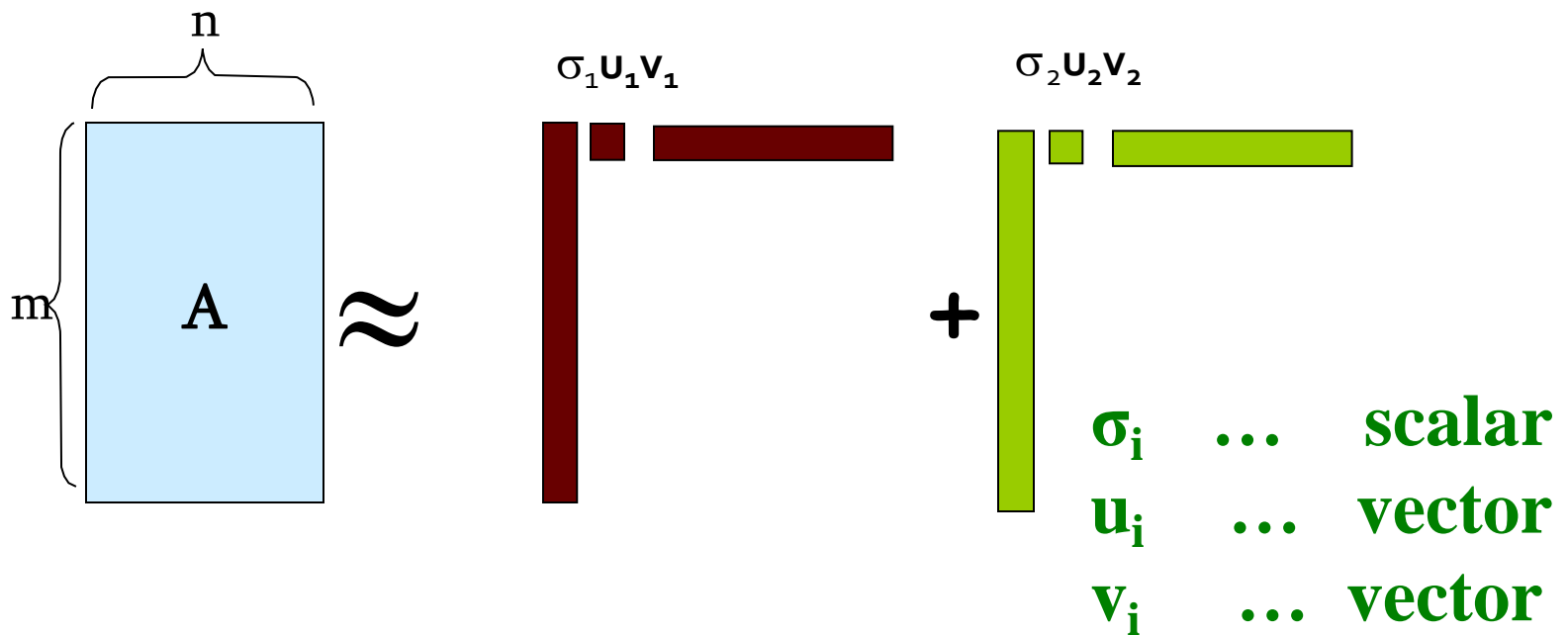
# SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^T$$



# SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^T$$



# SVD

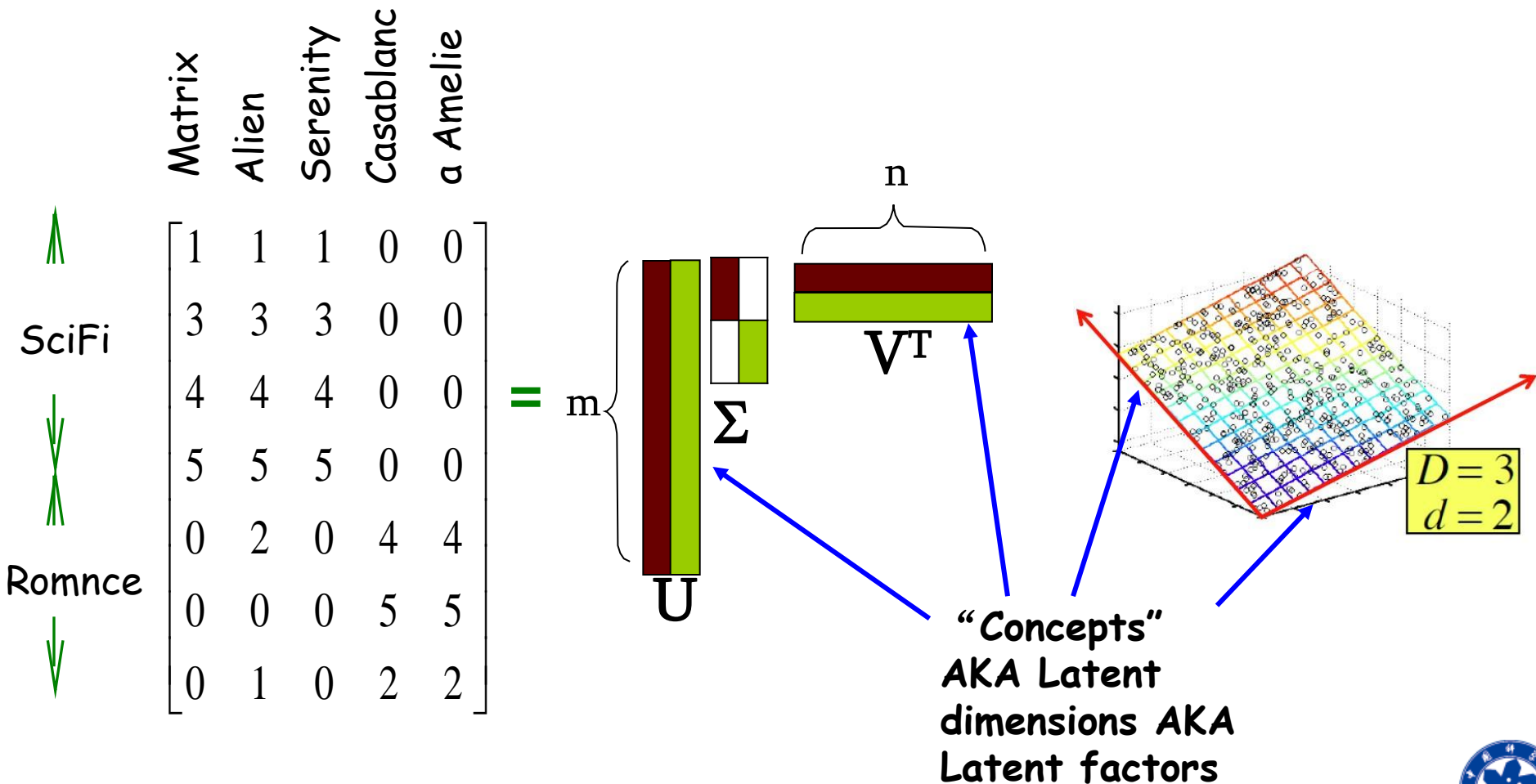
It is **always** possible to decompose a real matrix  $A$  into  $A = U\Sigma V^T$ , where

- $U, \Sigma, V$ : unique
- $U, V$ : column orthonormal
  - $U^T U = I; V^T V = I$  ( $I$ : identity matrix)
  - (Columns are orthogonal unit vectors)
- $\Sigma$ : diagonal
  - Entries (**singular values**) are positive,  
and sorted in decreasing order ( $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$ )



# SVD-Example:Users-to-Movies

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



# SVD-Example:Users-to-Movies

## ■ $A = U\Sigma V^T$ - example: Users to Movies

SciFi  $\uparrow$   
 $\downarrow$   
 Romnce  $\uparrow$   
 $\downarrow$

Matrix	Alien	Serenity	Casablanc	a 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

 $=$ 

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

 $\times$ 

12.4	0	0
0	9.5	0
0	0	1.3

 $\times$ 

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





# SVD-Example:Users-to-Movies

## ■ $A = U\Sigma V^T$ - example: Users to Movies

Matrix

	Alien	Serenity	Casablanca	Amelie
SciFi	1	1	0	0
	3	3	0	0
	4	4	0	0
	5	5	0	0
Romance	0	2	4	4
	0	0	5	5
	0	1	2	2

SciFi-concept

Romance-concept

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



# SVD-Example:Users-to-Movies

- $A = U\Sigma V^T$  - example:  $U$  is “user-to-concept” similarity matrix

Matrix

	Alien	Serenity	Casablanca	Amelie
SciFi	1	1	0	0
	3	3	0	0
	4	4	0	0
	5	5	0	0
Romnce	0	2	4	4
	0	0	5	5
	0	1	2	2

SciFi-concept Romance-concept

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



# SVD-Example:Users-to-Movies

■  $A = U\Sigma V^T$  - example:

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-concept

SciFi

Romance

SciFi-concept

"strength" of the SciFi-concept

12.4

9.5

1.3

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

0.56 0.59 0.56 0.09 0.09

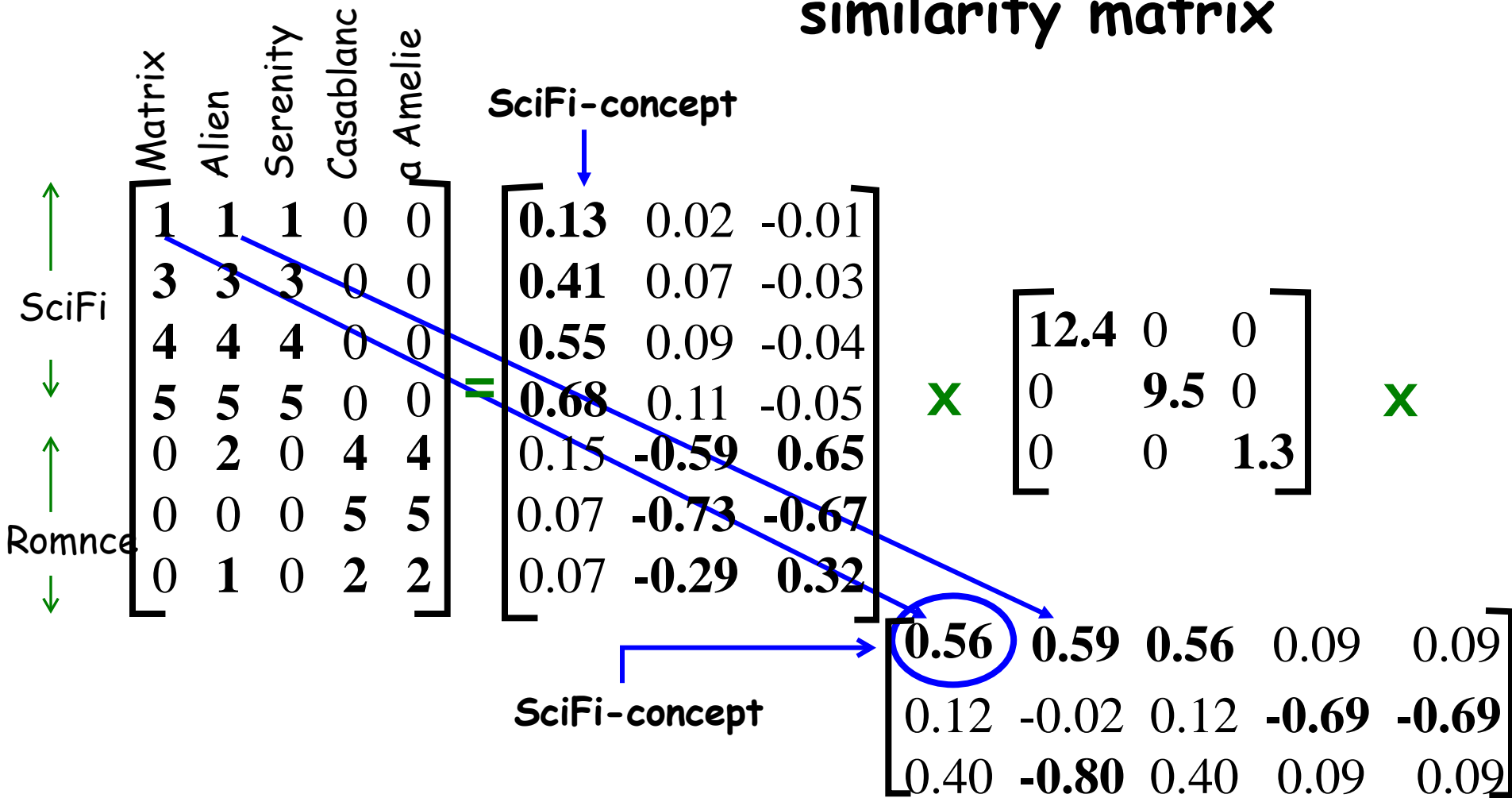
0.12 -0.02 0.12 -0.69 -0.69

0.40 -0.80 0.40 0.09 0.09



# SVD-Example:Users-to-Movies

- $A = U\Sigma V^T$  - example:  $V$  is "movie-to-concept" similarity matrix



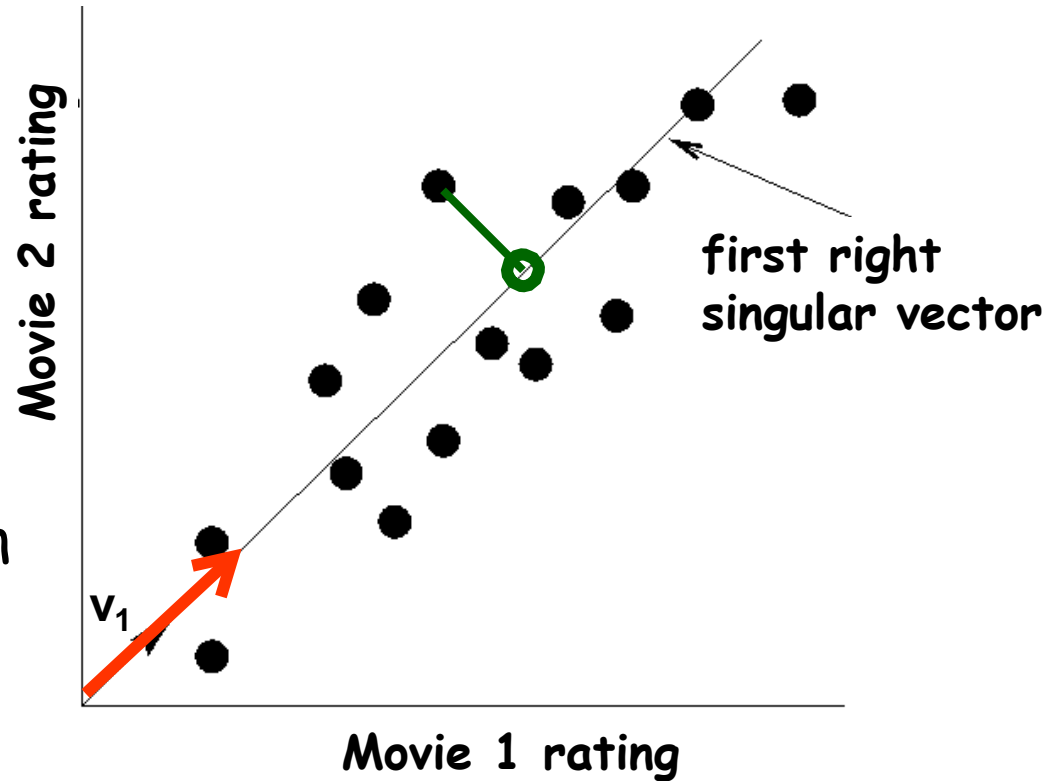
'movies', 'users' and 'concepts':

- $U$ : user-to-concept similarity matrix
- $V$ : movie-to-concept similarity matrix
- $\Sigma$ : its diagonal elements:  
'strength' of each concept



# SVD-Interpretation #2

- SVD gives 'best' axis to project on:
- 'best' = min sum of squares of projection errors
- In other words, minimum reconstruction error



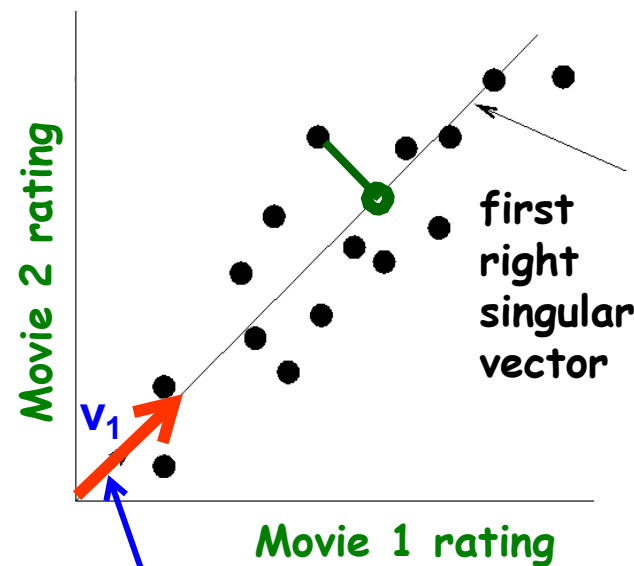
# SVD-Interpretation #2

- $A = U\Sigma V^T$  - example:
  - $V$ : "movie-to-concept" matrix
  - $U$ : "user-to-concept" matrix

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

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

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



# SVD-Interpretation #2

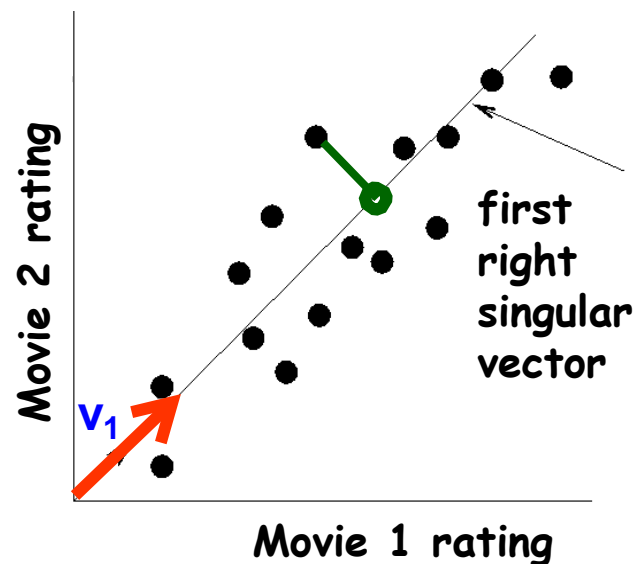
■  $A = U\Sigma V^T$  - example:

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

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

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

variance ('spread')  
on the  $v_1$  axis



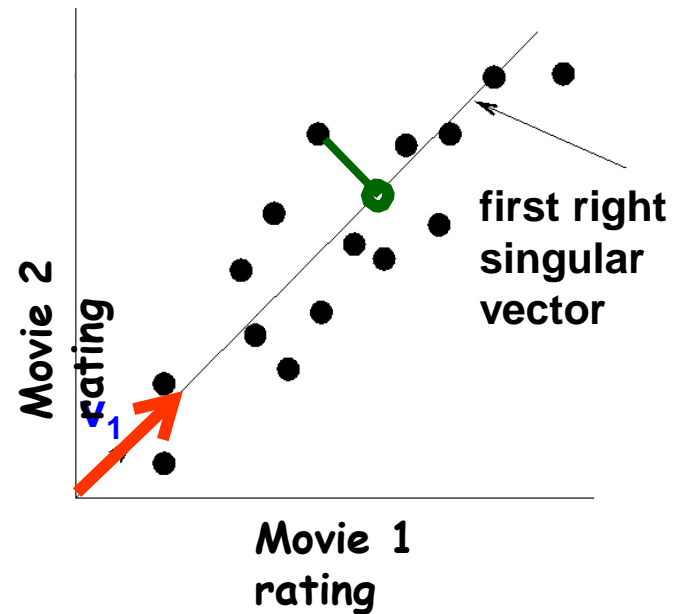


# SVD-Interpretation #2

- $A = U\Sigma V^T$  - example:
  - $U\Sigma$ : Gives the coordinates of the points in the projection axis

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix}$$

Projection of users  
on the "Sci-Fi"  
axis  $((U\Sigma)^T)$



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



# SVD-Interpretation #2

## More details

- Q: How exactly is dim. reduction done?

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



# SVD-Interpretation #2

## More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

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



# SVD-Interpretation #2

## More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

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



# SVD-Interpretation #2

## More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

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

The diagram illustrates the SVD decomposition of a matrix. The first matrix is a 7x5 matrix. The second matrix is a 7x3 matrix of singular values, with the third column (containing 1.3) crossed out with a red X. The third matrix is a 3x5 matrix of right singular vectors, with the third row (containing 0.40, -0.80, 0.40, 0.09, 0.09) crossed out with a red X. The green 'X' symbols indicate the multiplication of the matrices.



# SVD-Interpretation #2

## More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

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



# SVD-Interpretation #2

## More details

- **Q:** How exactly is dim. reduction done?
- **A:** Set smallest singular values to zero

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

Frobenius norm:

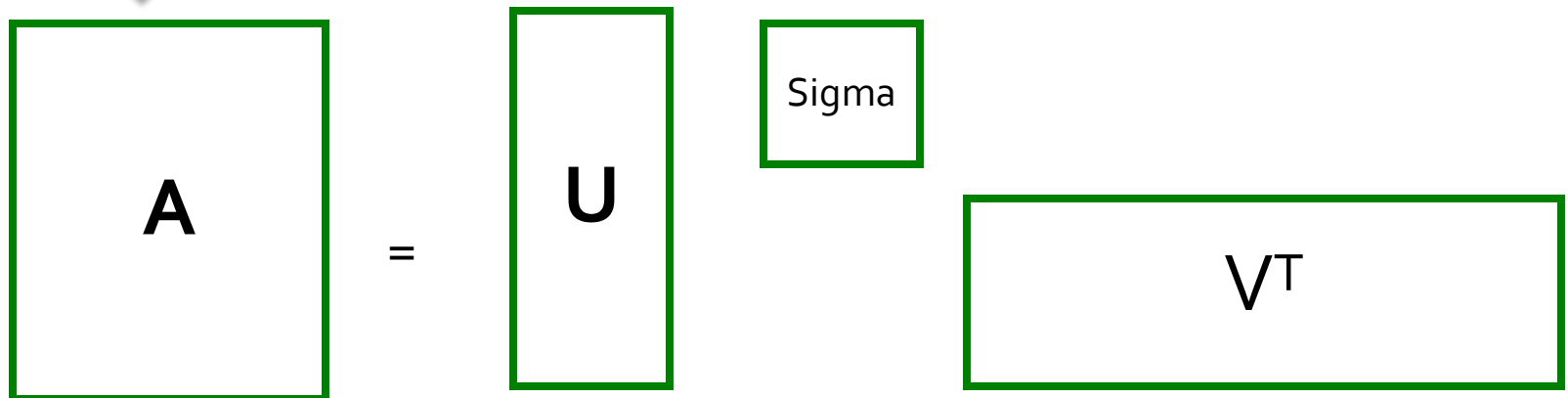
$$\|M\|_F = \sqrt{\sum_{ij} M_{ij}^2}$$

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

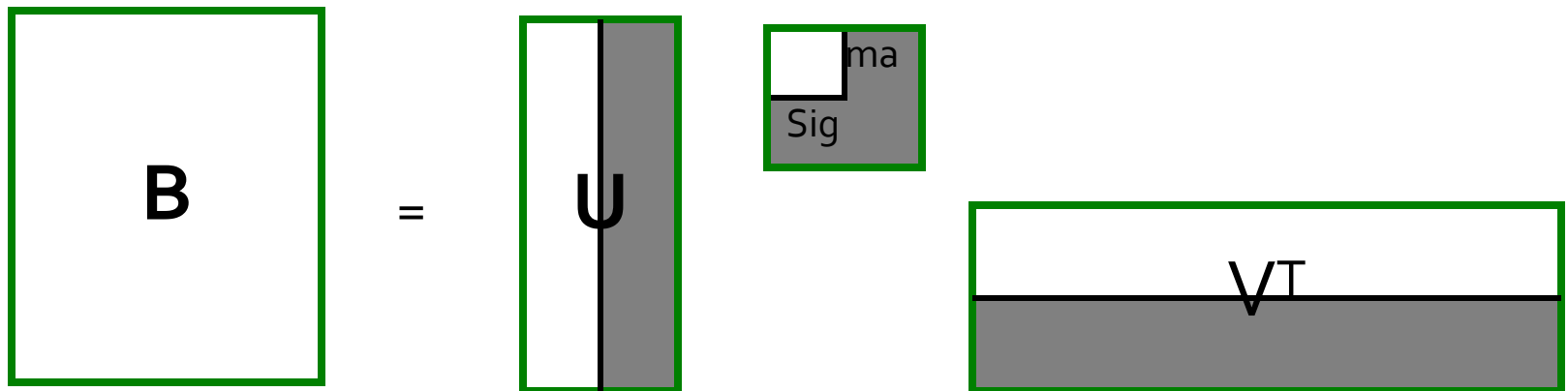
is “small”



# SVD-Best Low Rank Approx.



**B is best approximation of A**





# SVD-Best Low Rank Approx.

- Theorem: Let  $A = U \Sigma V^T$  ( $\sigma_1 \geq \sigma_2 \geq \dots$ ,  $\text{rank}(A)=r$ )  
then  $B = U S V^T$

- $S$  = diagonal  $n \times n$  matrix where  $s_i = \sigma_i$  ( $i=1 \dots k$ ) else  $s_i=0$   
**is a best rank- $k$  approximation to  $A$ :**

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

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & \\ \vdots & \vdots & \ddots & \\ x_{m1} & & & x_{mn} \end{pmatrix}_{m \times n} = \begin{pmatrix} u_{11} & \dots & & \\ \vdots & \ddots & & \\ u_{m1} & & & \end{pmatrix}_{m \times r} \begin{pmatrix} & & & \\ & 0 & \sigma_{11} & \\ & \ddots & & \\ \vdots & & & \end{pmatrix}_{r \times r} \begin{pmatrix} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \\ v_{r1} & \dots & \end{pmatrix}_{r \times n}$$

- **We will need 2 facts:**

- $\|M\|_F^2 = \sum (q_i)^2$  where  $M = PQR$  is SVD of  $M$
- $U \Sigma V^T - U S V^T = U (\Sigma - S) V^T$



# SVD-Best Low Rank Approx.

- **We will need 2 facts:**

- $\|M\|_F = \sqrt{\sum_k (q_k)^2}$  where  $M = P Q R$  is SVD of  $M$

$$\|M\| = \sum_i \sum_j (m_{ij})^2 = \sum_i \sum_j \left( \sum_k \sum_\ell p_{ik} q_{k\ell} r_{\ell j} \right)^2$$

$$\|M\| = \sum_i \sum_j \sum_k \sum_\ell \sum_n \sum_m p_{ik} q_{k\ell} r_{\ell j} p_{in} q_{nm} r_{mj}$$

$\sum_i p_{ik} p_{in}$  is 1 if  $k = n$  and 0 otherwise

**We apply:**

-- P column

orthonormal

-- R row orthonormal

-- Q is diagonal

- $U \Sigma V^T - U S V^T = U (\Sigma - S) V^T$



# SVD-Best Low Rank Approx.

- $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ ,  $\mathbf{B} = \mathbf{U} \mathbf{S} \mathbf{V}^T$  ( $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$ ,  $\text{rank}(\mathbf{A})=r$ )
    - $\mathbf{S}$  = diagonal  $n \times n$  matrix where  $s_i = \sigma_i$  ( $i=1 \dots k$ ) else  $s_i=0$
- then  $\mathbf{B}$  is solution to  $\min_{\mathbf{B}} \|\mathbf{A} - \mathbf{B}\|_F$ ,  $\text{rank}(\mathbf{B})=k$

## ■ Why?

$$\min_{\mathbf{B}, \text{rank}(\mathbf{B})=k} \|\mathbf{A} - \mathbf{B}\|_F = \min \|\mathbf{\Sigma} - \mathbf{S}\|_F = \min_s \sum_{i=1}^r (\sigma_i^2 - s_i^2)$$

We used:  $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T - \mathbf{U} \mathbf{S} \mathbf{V}^T = \mathbf{U} (\mathbf{\Sigma} - \mathbf{S}) \mathbf{V}^T$

- We want to choose  $s_i$  to minimize  $\sum (\sigma_i - s_i)^2$ 
  - We set  $s_i = \sigma_i$  ( $i=1 \dots k$ ) and other  $s_i=0$

$$= \min_{s_i} \sum_{i=1} (\sigma_i - s_i)^2 + \sum_{i=k+1} \sigma_i^2 = \sum_{i=k+1} \sigma_i^2$$



# SVD-Interpretation #2

Equivalent:

'spectral decomposition' of the matrix:

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



# SVD-Interpretation #2

Equivalent:

'spectral decomposition' of the matrix:

$$\begin{array}{c} \text{← m →} \\ \begin{array}{c} \text{↑} \\ \text{↓} \\ n \end{array} \left[ \begin{array}{ccccc} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{array} \right] \end{array} = \begin{array}{c} \text{← k terms →} \\ \sigma_1 \begin{array}{c} \text{u}_1 \\ \text{v}_1^T \end{array} + \sigma_2 \begin{array}{c} \text{u}_2 \\ \text{v}_2^T \end{array} + \dots \\ \begin{array}{cc} n \times 1 & 1 \times m \end{array} \end{array}$$

Assume:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$

Why is setting small  $\sigma_i$  to 0 the right thing to do?

Vectors  $u_i$  and  $v_i$  are unit length, so  $\sigma_i$  scales them.

So, zeroing small  $\sigma_i$  introduces less error



## SVD-Interpretation #2

**Q: How many  $\sigma_s$  to keep? A:**  
Rule-of-a thumb:

**keep 80-90% of 'energy' ( $=\sum \sigma_i^2$ )**

$$\begin{array}{c} \updownarrow \\ n \end{array} \begin{array}{c} \leftarrow m \rightarrow \\ \left| \begin{array}{ccccc} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{array} \right| \end{array} = \sigma_1 \mathbf{U}_1 \mathbf{V}_1^T + \sigma_2 \mathbf{U}_2 \mathbf{V}_2^T + \dots$$

Assume:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$



- **To compute SVD:**
  - $O(nm^2)$  or  $O(n^2m)$  (whichever is less)
- **But:**
  - Less work, if we just want singular values
  - or if we want first  $k$  singular vectors
  - or if the matrix is sparse
- **Implemented in linear algebra packages like**
  - LINPACK, Matlab, SPlus, Mathematica ...



# SVD-Conclusions so far

- **SVD:  $A = U \Sigma V^T$ : unique**
  - $U$ : user-to-concept similarities
  - $V$ : movie-to-concept similarities
  - $\Sigma$ : strength of each concept
- **Dimensionality reduction:**
  - keep the few largest singular values (80-90% of 'energy')
  - SVD: picks up linear correlations





# Relation to Eigen-decompositon

- **SVD gives us:**

- $A = U \Sigma V^T$

- **Eigen-decomposition:**

- $A = X \Lambda X^T$

- A is symmetric

- U, V, X are orthonormal ( $U^T U = I$ ),

- $\Lambda, \Sigma$  are diagonal

- **What is:**

- $AA^T = U \Sigma V^T (U \Sigma V^T)^T = U \Sigma V^T (V \Sigma^T \Sigma U^T) = U^T \Sigma \Sigma^T U^T$

- $A^T A = V \Sigma^T U^T (U \Sigma V^T) = V \Sigma \Sigma^T V^T$



# Relation to Eigen-decompositon

- **SVD gives us:**

- $A = U \Sigma V^T$

- **Eigen-decomposition:**

- $A = X \Lambda X^T$

- A is symmetric

- U, V, X are orthonormal ( $U^T U = I$ ),

- $\Lambda, \Sigma$  are diagonal

- **What is:**

- $AA^T = U \Sigma V^T (U \Sigma V^T)^T = U \Sigma V^T (V \Sigma^T U^T) = U \Sigma \Sigma^T U^T$

- $A^T A = V \Sigma^T U^T (U \Sigma V^T) = V \Sigma \Sigma^T V^T$

Shows how to  
compute SVD using  
eigenvalue  
decomposition!



$$\begin{matrix} X & \Lambda & X^T \\ \downarrow & \downarrow & \downarrow \end{matrix}$$

So,  $\lambda_i = \sigma_i^2$



# Case study: How to query?

		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				

$\begin{matrix} \uparrow \\ \text{SciFi} \\ \downarrow \\ \uparrow \\ \text{Romnce} \\ \downarrow \end{matrix}$

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

- Q: Find users that like 'Matrix'
- A: Map query into a 'concept space' - how?

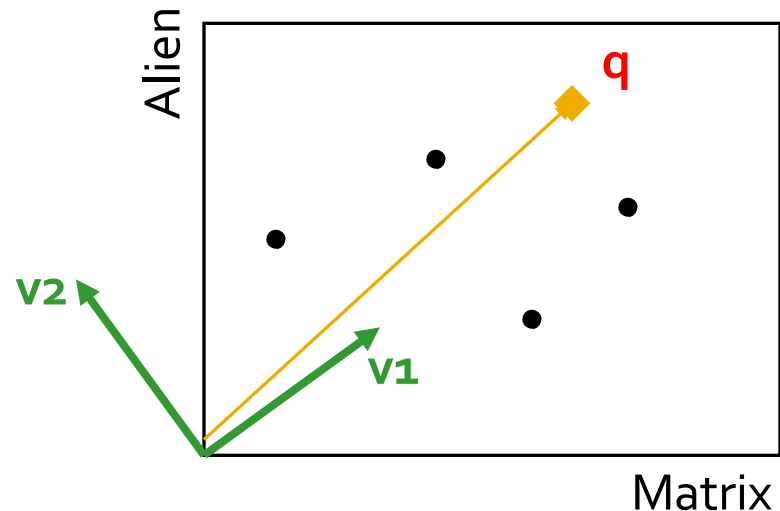


# Case study: How to query?

- Q: Find users that like 'Matrix'
- A: Map query into a 'concept space' – how?

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

**Project into concept space:**  
Inner product with each  
'concept' vector  $\mathbf{v}_i$

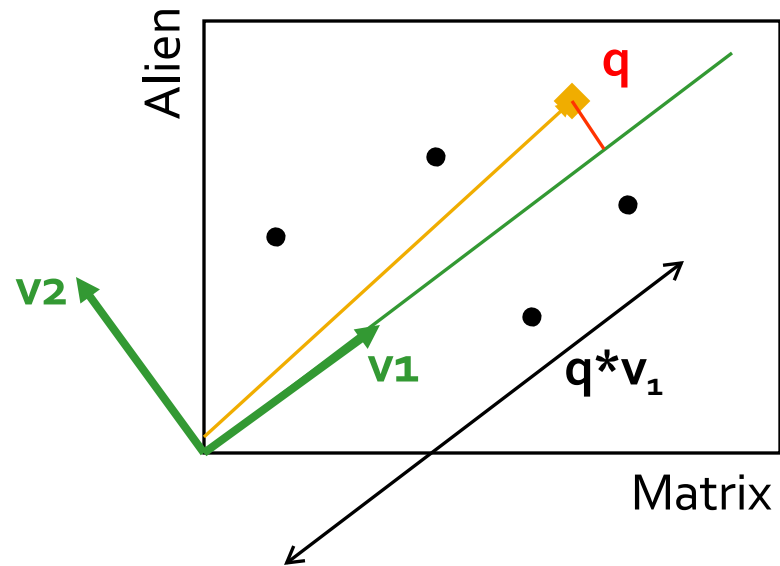


# Case study: How to query?

- Q: Find users that like 'Matrix'
- A: Map query into a 'concept space' – how?

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

**Project into concept space:**  
Inner product with each  
'concept' vector  $\mathbf{v}_i$



# Case study: How to query?

Compactly, we have:  $q_{\text{concept}} = q V$

E.g.:

$$q = \begin{bmatrix} \text{Matrix} \\ \text{Alien} \\ \text{Serenity} \\ \text{Casablanca} \\ \text{Amelie} \end{bmatrix} \begin{matrix} 5 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \times \begin{bmatrix} 0.56 & 0.12 \\ 0.59 & -0.02 \\ 0.56 & 0.12 \\ 0.09 & -0.69 \\ 0.09 & -0.69 \end{bmatrix} = \begin{bmatrix} \text{SciFi-concept} \\ 2.8 \end{bmatrix} \begin{matrix} 0.6 \end{matrix}$$

movie-to-concept similarities (V)



# Case study: How to query?

- How would the user  $d$  that rated ('Alien', 'Serenity') be handled?

$$\mathbf{d}_{\text{concept}} = \mathbf{d} \mathbf{V}$$

E.g.:

$$\mathbf{q} = \begin{matrix} & \text{Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ \begin{bmatrix} 0 & 4 & 5 & 0 & 0 \end{bmatrix} & \times & \begin{bmatrix} 0.56 & 0.12 \\ 0.59 & -0.02 \\ 0.56 & 0.12 \\ 0.09 & -0.69 \\ 0.09 & -0.69 \end{bmatrix} & = & \begin{bmatrix} 5.2 & 0.4 \end{bmatrix} \end{matrix}$$

movie-to-concept  
similarities (V)

SciFi-concept  
↓



# Case study: How to query?

- **Observation:** User  $d$  that rated ('*Alien*', '*Serenity*') will be **similar** to user  $q$  that rated ('*Matrix*'), although  $d$  and  $q$  have **zero ratings in common**!

	Matrix	Alien	Serenity	Casablanca	Amelie	
$d =$	0	4	5	0	0	$\longrightarrow$ $\begin{bmatrix} 2.8 & 0.6 \end{bmatrix}$
$q =$	5	0	0	0	0	$\longrightarrow$ $\begin{bmatrix} 5.2 & 0.4 \end{bmatrix}$
Zero ratings in common						Similarity $\neq 0$

SciFi-concept





## SVD: Drawbacks

- + Optimal low-rank approximation

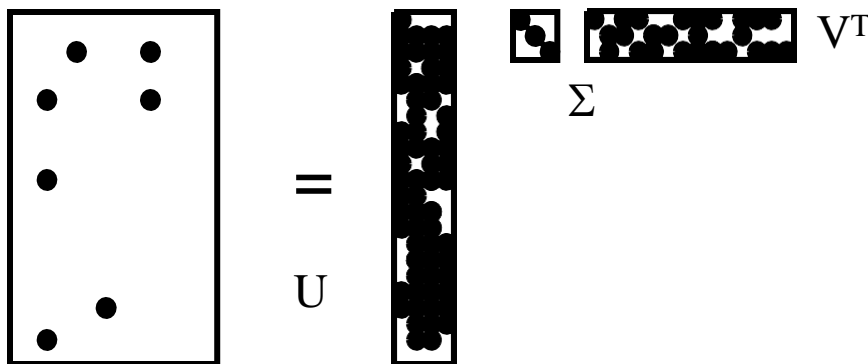
in terms of Frobenius norm

- Interpretability problem:

- A singular vector specifies a linear combination of all input columns or rows

- Lack of sparsity:

- Singular vectors are **dense**!



## Announcements:

- HW2 has been posted
- LSH Gradiance quiz has been posted. Due 2013-01-30 23:59
- Date for an alternate final: Tue 3/19 6-9PM

## CUR Decomposition



# CUR Decomposition

Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- Goal: Express  $A$  as a product of matrices  $C, U, R$  Make  $\|A - C \cdot U \cdot R\|_F$  small
- "Constraints" on  $C$  and  $R$ :

$$\left( \begin{array}{|c|} \hline \text{Red} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{Blue} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{Dark Red} \\ \hline \end{array} \right) \approx \left( \begin{array}{|c|c|c|c|c|c|} \hline \text{Red} & \text{Red} & \text{Red} & \text{Blue} & \text{Dark Red} & \text{Dark Red} \\ \hline \end{array} \right) \cdot \left( \begin{array}{c} U \end{array} \right) \cdot \left( \begin{array}{c} R \end{array} \right)$$

$A$ 
 $C$ 
 $U$ 
 $R$

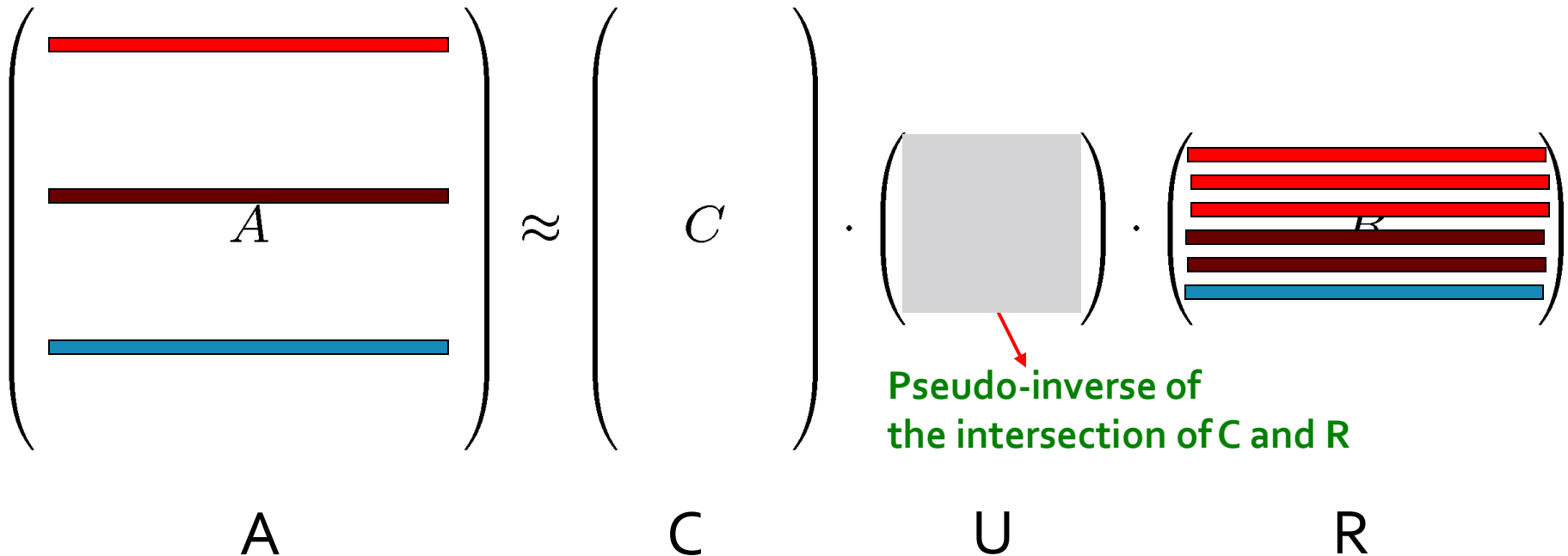


# CUR Decomposition

Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- Goal: Express  $A$  as a product of matrices  $C, U, R$  Make  $\|A - C \cdot U \cdot R\|_F$  small
- "Constraints" on  $C$  and  $R$ :



Pseudo-inverse of  
the intersection of C and R



# CUR: Provably good approx.to SVD

- **Let:**

$A_k$  be the "best" rank  $k$  approximation to  $A$  (that is,  $A_k$  is SVD of  $A$ )

## Theorem [Drineas et al.]

**CUR** in  $O(m \cdot n)$  time achieves

- $\|A - CUR\|_F \leq \|A - A_k\|_F + \delta \|A\|_F$

with probability at least  $1 - \delta$  by picking

- $O(k \log(1/\delta)/\epsilon)$  columns, and
- $O(k^2 \log^3(1/\delta)/\epsilon)$  rows

**In practice:**

Pick  $4k$   
cols/rows



# CUR: How it Works

## ■ Sampling columns (similarly for rows):

**Input:** matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , sample size  $c$

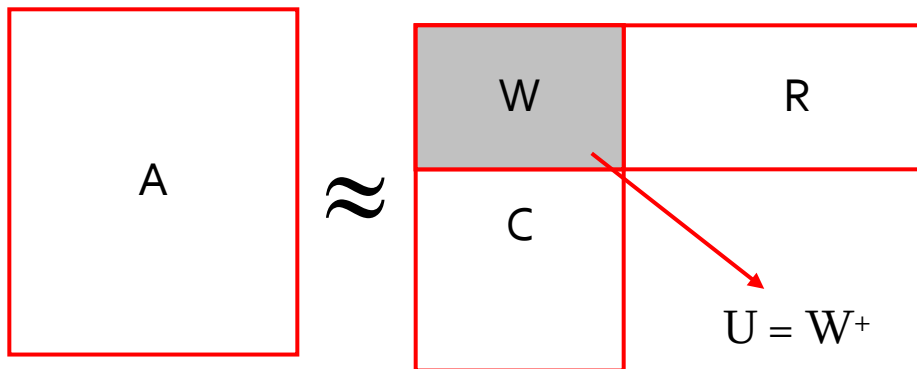
**Output:**  $\mathbf{C}_d \in \mathbb{R}^{m \times c}$

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  $\mathbf{C}_d(:, i) = \mathbf{A}(:, j) / \sqrt{cP(j)}$



# Computing U

- Let **W** be the “intersection” of sampled columns **C** and rows **R**
  - Let SVD of  $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}^T$
- **Then:  $\mathbf{U} = \mathbf{W}^+ = \mathbf{Y} \mathbf{Z}^+ \mathbf{X}^T$** 
  - $\mathbf{Z}^+$ : reciprocals of non-zero singular values:  $Z^+_{ii} = 1/Z_{ii}$
  - $\mathbf{W}^+$  is the “pseudoinverse”



## Why pseudoinverse works?

$\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}$  then  $\mathbf{W}^{-1} = \mathbf{X}^{-1} \mathbf{Z}^{-1} \mathbf{Y}^{-1}$

Due to orthonormality  $\mathbf{X}^{-1} = \mathbf{X}^T$  and  $\mathbf{Y}^{-1} = \mathbf{Y}^T$

Since **Z** is diagonal  $Z^{-1} = 1/Z_{ii}$

**Thus**, if **W** is nonsingular, pseudoinverse is the true inverse



# CUR: Pros&Cons

+ Easy interpretation

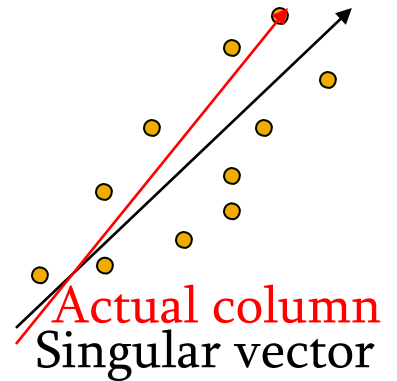
- Since the basis vectors are actual columns and rows

+ **Sparse basis**

- Since the basis vectors are actual columns and rows

- **Duplicate columns and rows**

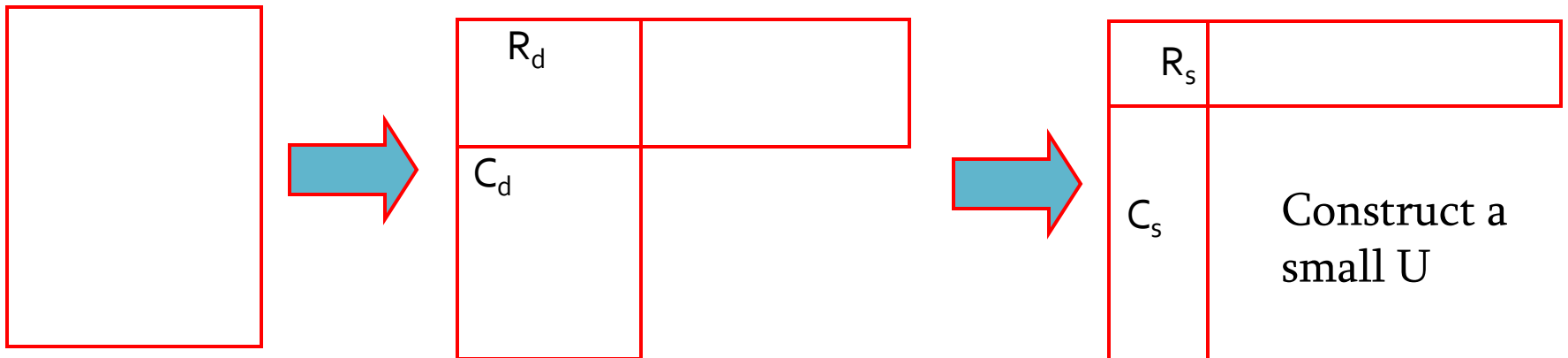
- Columns of large norms will be sampled many times





# Solution

- If we want to get rid of the duplicates:
  - Throw them away
  - Scale (multiply) the columns/rows by the square root of the number of duplicates



# SVD VS. CUR

SVD:  $A = U \Sigma V^T$

Annotations for SVD:

- $A$ : Huge but sparse
- $U$ : Big and dense
- $\Sigma$ : sparse and small
- $V^T$ : Big and dense

CUR:  $A = C U R$

Annotations for CUR:

- $A$ : Huge but sparse
- $C$ : Big but sparse
- $U$ : dense but small
- $R$ : Big but sparse



# Simple Experiment

## ■ DBLP bibliographic data

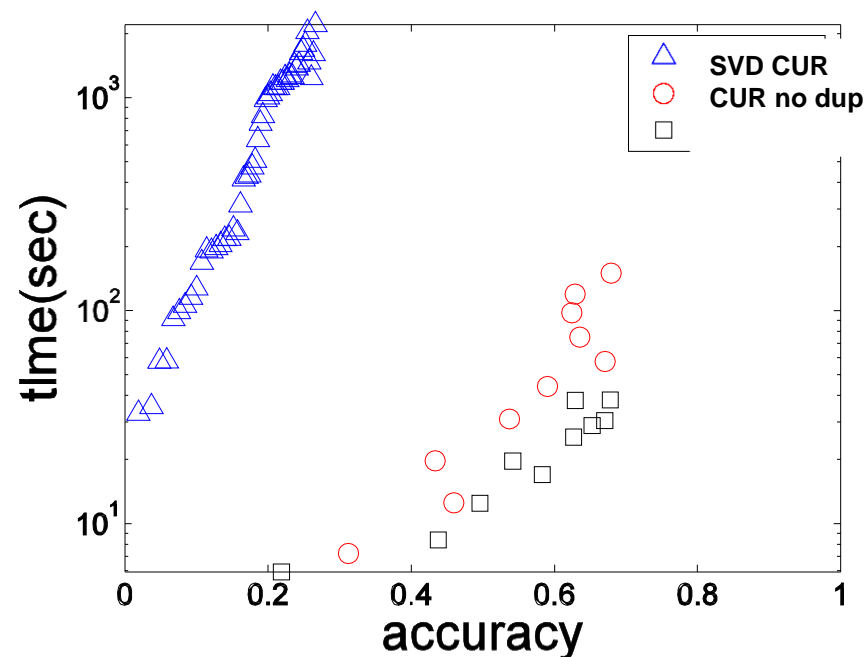
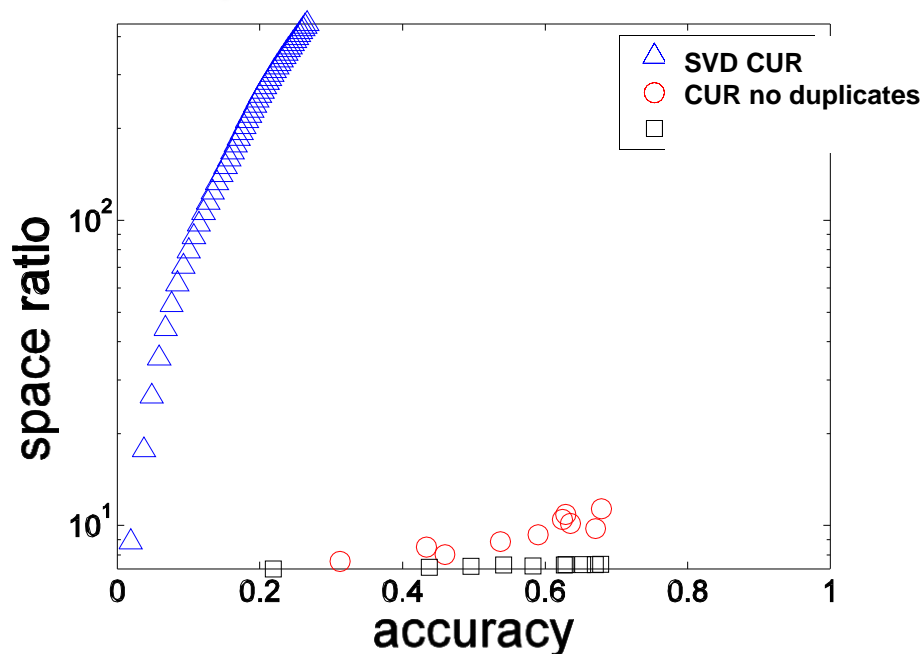
- Author-to-conference big sparse matrix
- $A_{ij}$ : Number of papers published by author  $i$  at conference  $j$
- 428K authors (rows), 3659 conferences (columns)
  - Very sparse

## ■ Want to reduce dimensionality

- How much time does it take?
- What is the reconstruction error?
- How much space do we need?



# Results : DBLP-big sparse matrix



- **Accuracy:**
  - 1 - relative sum squared errors
- **Space ratio:**
  - $\# \text{output matrix entries} / \# \text{input matrix entries}$
- **CPU time**



# What about linearity assumption?

- **SVD is limited to linear projections:**

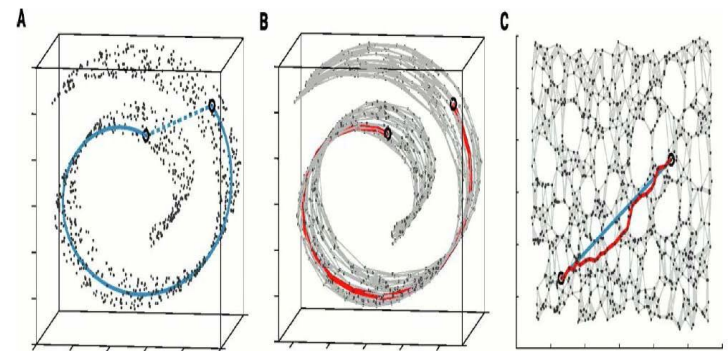
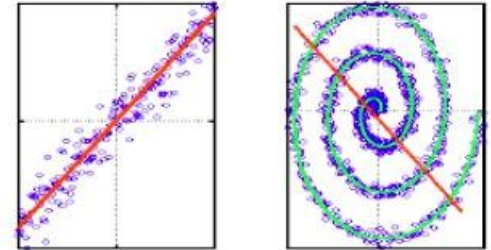
- Lower-dimensional linear projection that preserves Euclidean distances

- Non-linear methods: **Isomap**

- Data lies on a nonlinear low-dim curve aka manifold
  - Use the distance as measured along the manifold

- **How?**

- Build adjacency graph
- Geodesic distance is graph distance
- SVD/PCA the graph pairwise distance matrix



## Further Reading: CUR

- Drineas et al., *Fast Monte Carlo Algorithms for Matrices III: Computing a Compressed Approximate Matrix Decomposition*, SIAM Journal on Computing, 2006.
- J. Sun, Y. Xie, H. Zhang, C. Faloutsos: *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM 2007
- *Intra- and interpopulation genotype reconstruction from tagging SNPs*, P. Paschou, M. W. Mahoney, A. Javed, J. R. Kidd, A. J. Pakstis, S. Gu, K. K. Kidd, and P. Drineas, Genome Research, 17(1), 96-107 (2007)
- *Tensor-CUR Decompositions For Tensor-Based Data*, M. W. Mahoney, M. Maggioni, and P. Drineas, Proc. 12-th Annual SIGKDD, 327-336 (2006)



Tell me and I forget.  
Show me and I remember.  
Involve me and I understand.

Thank you! Q&A

