

# DATA MINING

# LECTURE 7

---

Dimensionality Reduction  
PCA – SVD

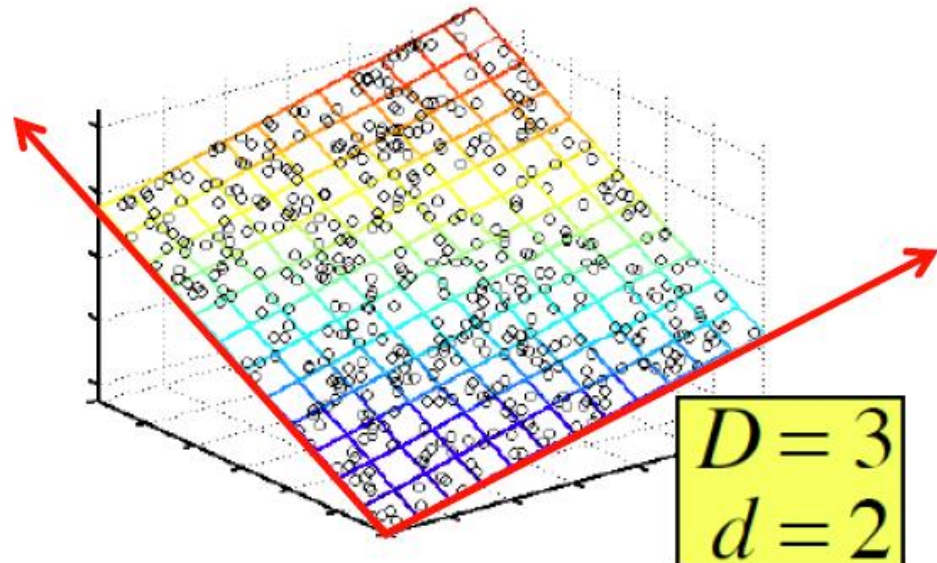
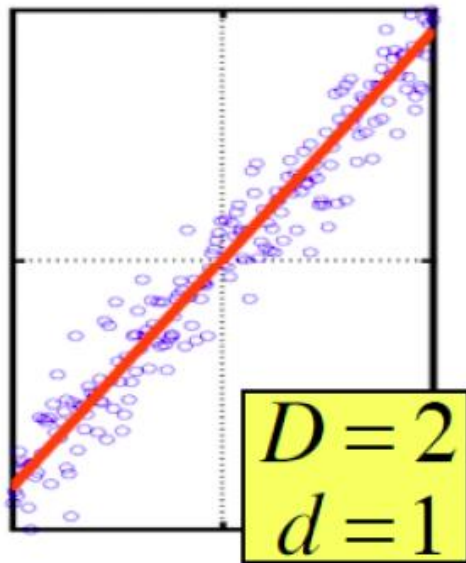
(Thanks to Jure Leskovec, Evimaria Terzi)

# The curse of dimensionality

- Real data usually have **thousands**, or **millions** of dimensions
  - E.g., web documents, where the dimensionality is the vocabulary of words
  - Facebook graph, where the dimensionality is the number of users
- Huge number of dimensions causes problems
  - Data becomes very **sparse**, some algorithms become meaningless (e.g. density based clustering)
  - The **complexity** of several algorithms depends on the dimensionality and they become infeasible (e.g. nearest neighbor search).

# Dimensionality Reduction

- Usually the data can be described with fewer dimensions, without losing much of the meaning of the data.
  - The data **reside** in a space of lower dimensionality



# Example

TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

- In this data matrix the dimension is essentially **3**
  - There are **three** types of products and **three** types of users

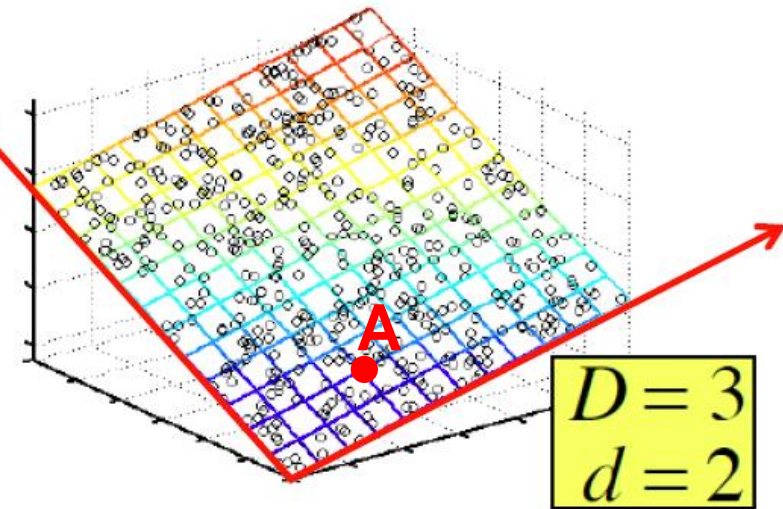
# Example

- **Cloud of points 3D space:**

- Think of point positions

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

1 row per point:



- **We can rewrite coordinates more efficiently!**

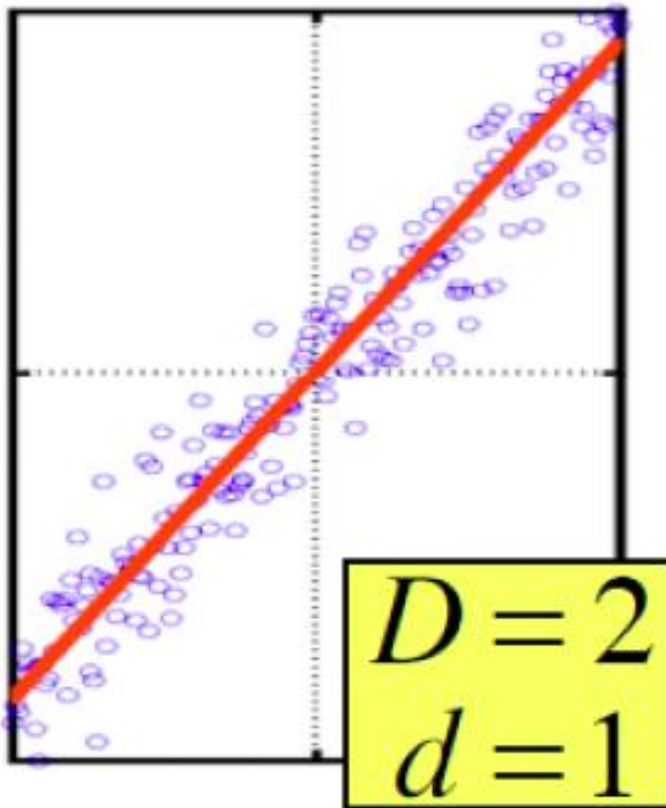
- Old basis vectors:  $[1 \ 0 \ 0]$   $[0 \ 1 \ 0]$   $[0 \ 0 \ 1]$
- **New basis vectors:  $[1 \ 2 \ 1]$   $[-2 \ -3 \ 1]$**
- Then **A** has new coordinates:  $[1 \ 0]$ . **B**:  $[0 \ 1]$ , **C**:  $[1 \ -1]$ 
  - **Notice: We reduced the number of coordinates!**

# Dimensionality Reduction

- Find the “**true dimension**” of the data
  - In reality things are never as clear and simple as in this example, but we can still reduce the dimension.
- Essentially, we assume that some of the data is **noise**, and we can approximate the useful part with a lower dimensionality space.
  - Dimensionality reduction does not just reduce the amount of data, it often brings out the **useful** part of the data

# Dimensionality Reduction

- **Goal of dimensionality reduction is to discover the axis of data!**

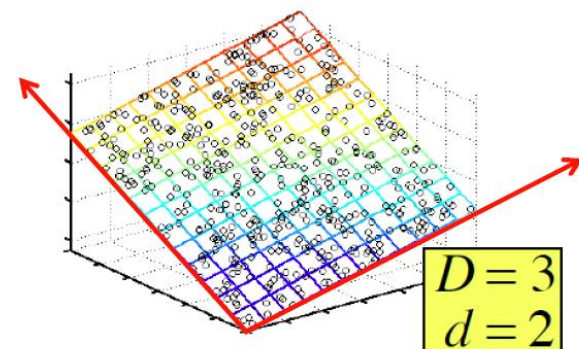


Rather than representing every point with 2 coordinates we represent each point with 1 coordinate (corresponding to the position of the point on the red line).

By doing this we incur a bit of **error** as the points do not exactly lie on the line

# Why Reduce Dimensions?

- **Discover hidden correlations/topics**
  - E.g., in documents, words that occur commonly together
- **Remove redundant and noisy features**
  - E.g., in documents, not all words are useful
- **Interpretation and visualization**
- **Easier storage and processing of the data**





# Dimensionality Reduction

- We have already seen a form of dimensionality reduction
- LSH, and random projections reduce the dimension while preserving the distances

# Data in the form of a matrix

- We are given  $n$  objects and  $d$  attributes describing the objects. Each object has  $d$  numeric values describing it.
- We will represent the data as a  $n \times d$  real matrix  $A$ .
  - We can now use tools from **linear algebra** to process the data matrix
- Our goal is to produce a new  $n \times k$  matrix  $B$  such that
  - It preserves as much of the **information** in the original matrix  $A$  as possible
  - It reveals something about the structure of the data in  $A$

# Example: Document matrices

**d** terms

(e.g., theorem, proof, etc.)

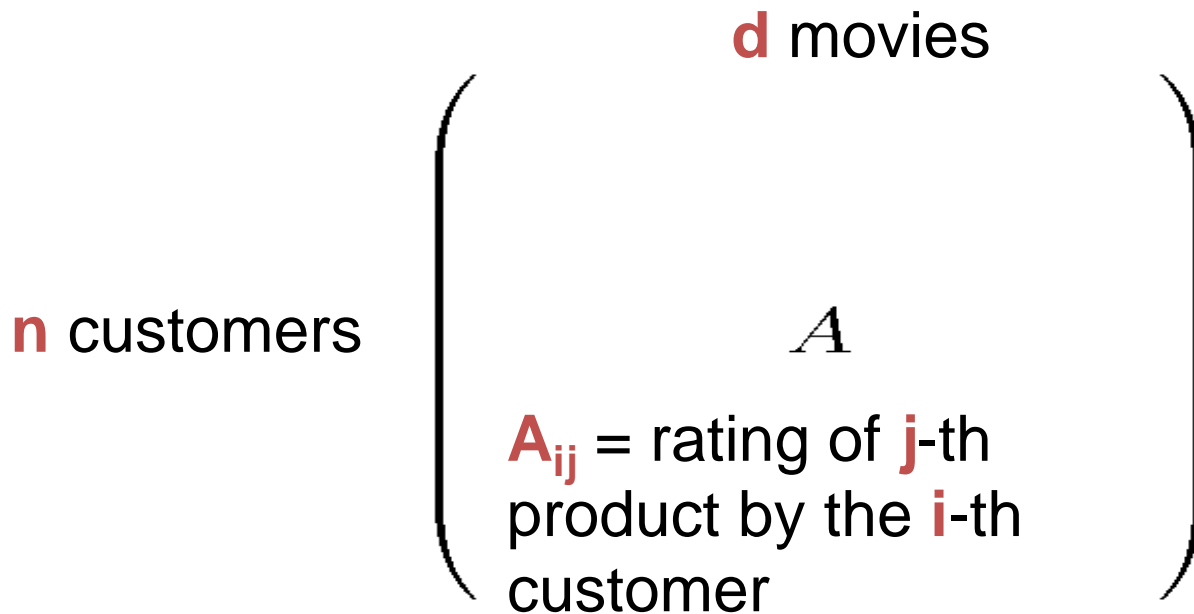
**n** documents

$A$

$A_{ij}$  = frequency of the **j**-th term in the **i**-th document

Find subsets of terms that bring documents together

# Example: Recommendation systems



Find subsets of movies that capture the behavior of the customers

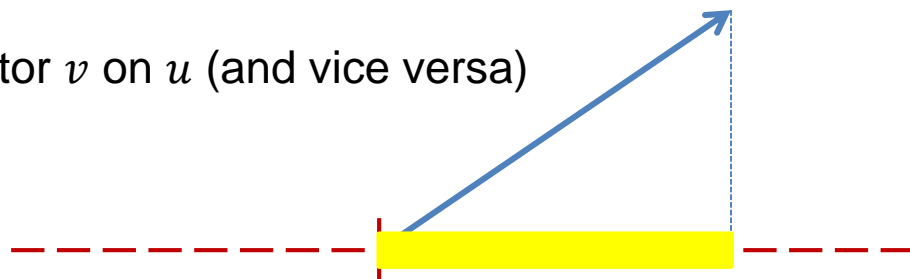
# Linear algebra

- We assume that vectors are **column vectors**.
- We use  $v^T$  for the **transpose** of vector  $v$  (**row vector**)

- **Dot product**:  $u^T v$  ( $1 \times n, n \times 1 \rightarrow 1 \times 1$ )

- The dot product is the **projection** of vector  $v$  on  $u$  (and vice versa)

- $[1, 2, 3] \begin{bmatrix} 4 \\ 1 \\ 2 \end{bmatrix} = 12$



- $u^T v = \|v\| \|u\| \cos(u, v)$

- If  $\|u\| = 1$  (**unit vector**) then  $u^T v$  is the **projection length** of  $v$  on  $u$

- $[-1, 2, 3] \begin{bmatrix} 4 \\ -1 \\ 2 \end{bmatrix} = 0$  : **orthogonal** vectors

- **Orthonormal** vectors: two unit vectors that are orthogonal

# Matrices

- An  $n \times m$  matrix  $A$  is a collection of  $n$  row vectors and  $m$  column vectors

$$A = \begin{bmatrix} | & | & | \\ a_1 & a_2 & a_3 \\ | & | & | \end{bmatrix} \quad A = \begin{bmatrix} - & \alpha_1^T & - \\ - & \alpha_2^T & - \\ - & \alpha_3^T & - \end{bmatrix}$$

- Matrix-vector multiplication
  - **Right multiplication**  $Au$ : **projection** of  $u$  onto the **row vectors** of  $A$ , or projection of row vectors of  $A$  onto  $u$ .
  - **Left-multiplication**  $u^T A$ : **projection** of  $u$  onto the **column vectors** of  $A$ , or projection of column vectors of  $A$  onto  $u$
- Example:

$$[1,2,3] \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = [1,2]$$

# Rank

- **Row space** of A: The set of vectors that can be written as a linear combination of the **rows** of A
  - All vectors of the form  $v = u^T A$
- **Column space** of A: The set of vectors that can be written as a linear combination of the **columns** of A
  - All vectors of the form  $v = Au$ .
- **Rank** of A: the number of **linearly independent** row (or column) vectors
  - These vectors define a **basis** for the row (or column) space of A
    - All vectors in the row (column) space are linear combinations of the basis vectors
- **Example**
  - Matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$  has rank  $\mathbf{r=2}$ 
    - **Why?** The first two rows are linearly independent, so the rank is at least 2, but all three rows are linearly dependent (the first is equal to the sum of the second and third) so the rank must be less than 3.

# Rank-1 matrices

- In a rank-1 matrix, all columns (or rows) are multiples of the same column (or row) vector

$$A = \begin{bmatrix} 1 & 2 & -1 \\ 2 & 4 & -2 \\ 3 & 6 & -3 \end{bmatrix}$$

- All **rows** are multiples of  $r^T = [1, 2, -1]$

- All **columns** are multiples of  $c = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

- **External product:**  $uv^T$  ( $n \times 1, 1 \times m \rightarrow n \times m$ )

- The resulting  $n \times m$  has **rank 1**: all rows (or columns) are **linearly dependent**

- $A = cr^T$



# Eigenvectors

- (Right) Eigenvector of matrix  $A$ : a vector  $v$  such that  $Av = \lambda v$
- $\lambda$ : eigenvalue of eigenvector  $v$
- A square symmetric matrix  $A$  of rank  $r$ , has  $r$  orthonormal eigenvectors  $u_1, u_2, \dots, u_r$  with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_r$ .
- Eigenvectors define an orthonormal basis for the column space of  $A$
- We can write:

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

# Singular Value Decomposition

$$A = U \Sigma V^T = [u_1, u_2, \dots, u_r] \begin{bmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_r \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_r^T \end{bmatrix}$$

$[n \times m] = [n \times r] [r \times r] [r \times m]$   
 $r$ : rank of matrix  $A$

- $\sigma_1, \geq \sigma_2 \geq \dots \geq \sigma_r$ : singular values of matrix  $A$
- $u_1, u_2, \dots, u_r$ : left singular vectors of  $A$
- $v_1, v_2, \dots, v_r$ : right singular vectors of  $A$

$$A = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_r u_r v_r^T$$

# Singular Value Decomposition

- The left singular vectors are an orthonormal basis for the row space of  $A$ .
- The right singular vectors are an orthonormal basis for the column space of  $A$ .
- If  $A$  has rank  $r$ , then  $A$  can be written as the sum of  $r$  rank-1 matrices
- There are  $r$  “linear components” (trends) in  $A$ .
  - Linear trend: the tendency of the row vectors of  $A$  to align with vector  $\mathbf{v}$
  - Strength of the  $i$ -th linear trend:  $\|A\mathbf{v}_i\| = \sigma_i$

# Symmetric matrices

- Special case:  $A$  is **symmetric** positive definite matrix

$$A = \lambda_1 u_1 u_1^T + \lambda_2 u_2 u_2^T + \cdots + \lambda_r u_r u_r^T$$

- $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_r \geq 0$ : **Eigenvalues** of  $A$
- $u_1, u_2, \dots, u_r$ : **Eigenvectors** of  $A$

# An (extreme) example

- User-Movie matrix
  - Blue and Red rows (columns) are **linearly dependent**

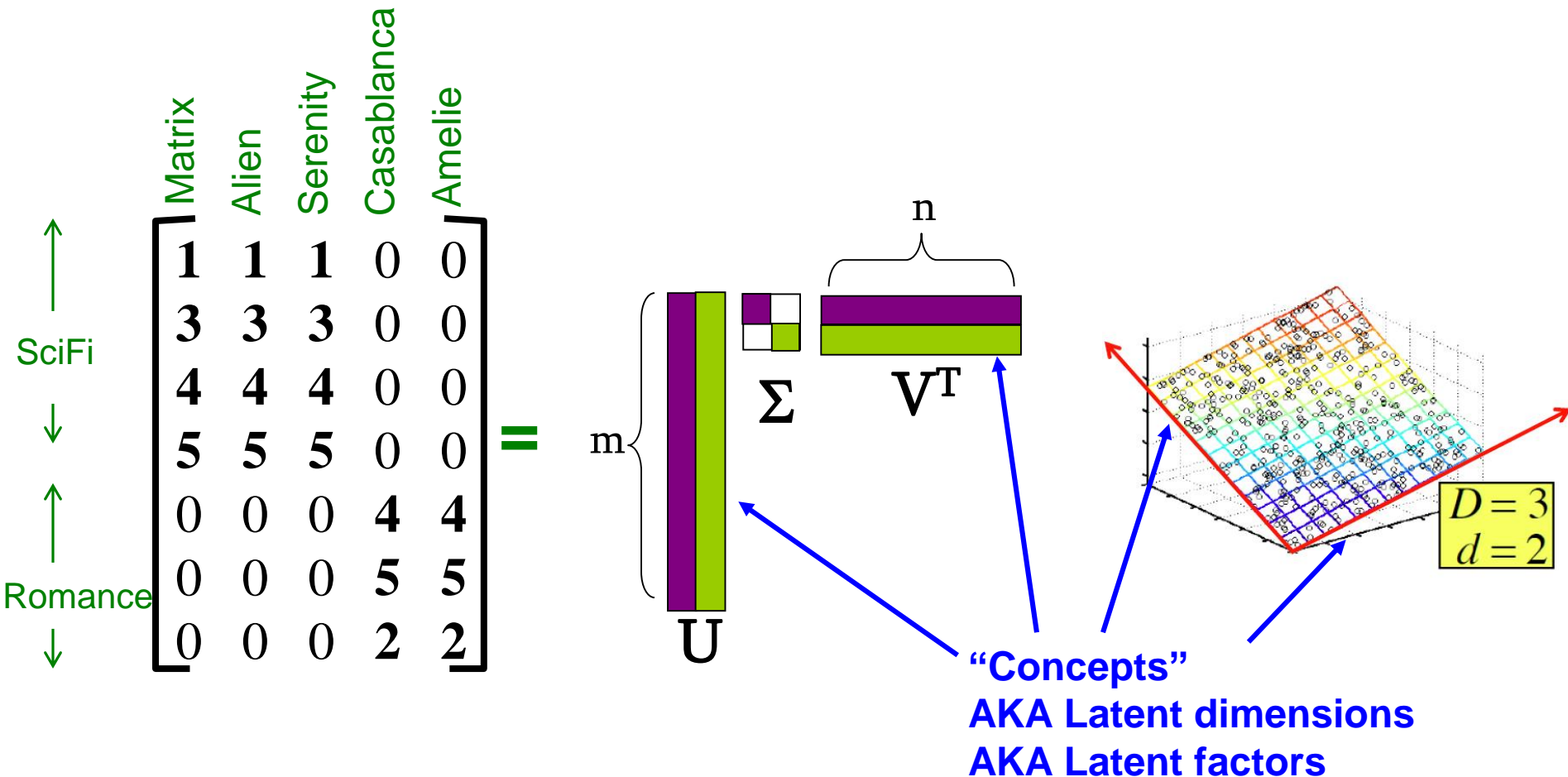
$$A = \begin{array}{|c|c|} \hline \text{Blue} & \text{White} \\ \hline \text{White} & \text{Red} \\ \hline \end{array}$$

- There are two **prototype** users (vectors of movies): blue and red
  - To describe the data is enough to describe the two **prototypes**, and the **projection weights** for each row
- $A$  is a **rank-2** matrix

$$A = [w_1, w_2] \begin{bmatrix} d_1^T \\ d_2^T \end{bmatrix}$$

# 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

The diagram illustrates the SVD decomposition of a user-movie rating matrix  $A$  into three matrices:  $U$ ,  $\Sigma$ , and  $V^T$ .

**Matrix A (User-Movie Ratings):**

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

**Matrix  $\Sigma$  (Concepts):**

SciFi-concept	0.14	0.00
	0.42	0.00
	0.56	0.00
	0.70	0.00
	0.00	0.60
	0.00	0.75
	0.00	0.30

**Matrix  $V^T$  (Concepts):**

	12.4	0
	0	9.5

**Matrix  $U$  (User Latent Space):**

	0.58	0.58	0.58	0.00	0.00
	0.00	0.00	0.00	0.71	0.71

The decomposition is shown as:  $A = U \Sigma V^T$ .

# SVD – Example: Users-to-Movies

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

Matrix A (User-to-Movie Ratings):

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

Matrix U (User-to-Concept Similarity):

	SciFi	Romance
0.14	0.00	
0.42	0.00	
0.56	0.00	
0.70	0.00	
0.00	0.60	
0.00	0.75	
0.00	0.30	

Matrix  $\Sigma$  (Singular Values):

12.4	0
0	9.5

Matrix  $V^T$  (Concept-to-Movie Similarity):

0.58	0.58	0.58	0.00	0.00
0.00	0.00	0.00	0.71	0.71

$U$  is "user-to-concept" similarity matrix



# SVD – Example: Users-to-Movies

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

Matrix Alien Serenity Casablanca Amelie

SciFi ↑

↓

Romance ↑

↓

SciFi-concept

Romance-concept

$V$  is “movie to concept” similarity matrix

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.14 & 0.00 \\ 0.42 & 0.00 \\ 0.56 & 0.00 \\ 0.70 & 0.00 \\ 0.00 & 0.60 \\ 0.00 & 0.75 \\ 0.00 & 0.30 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.71 & 0.71 \end{bmatrix}$$

# SVD – Example: Users-to-Movies

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

SciFi ↑

Romance ↓

	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	0	0	4	4	
	0	0	0	5	5	
↓	0	0	0	2	2	

SciFi-concept

Romance-concept

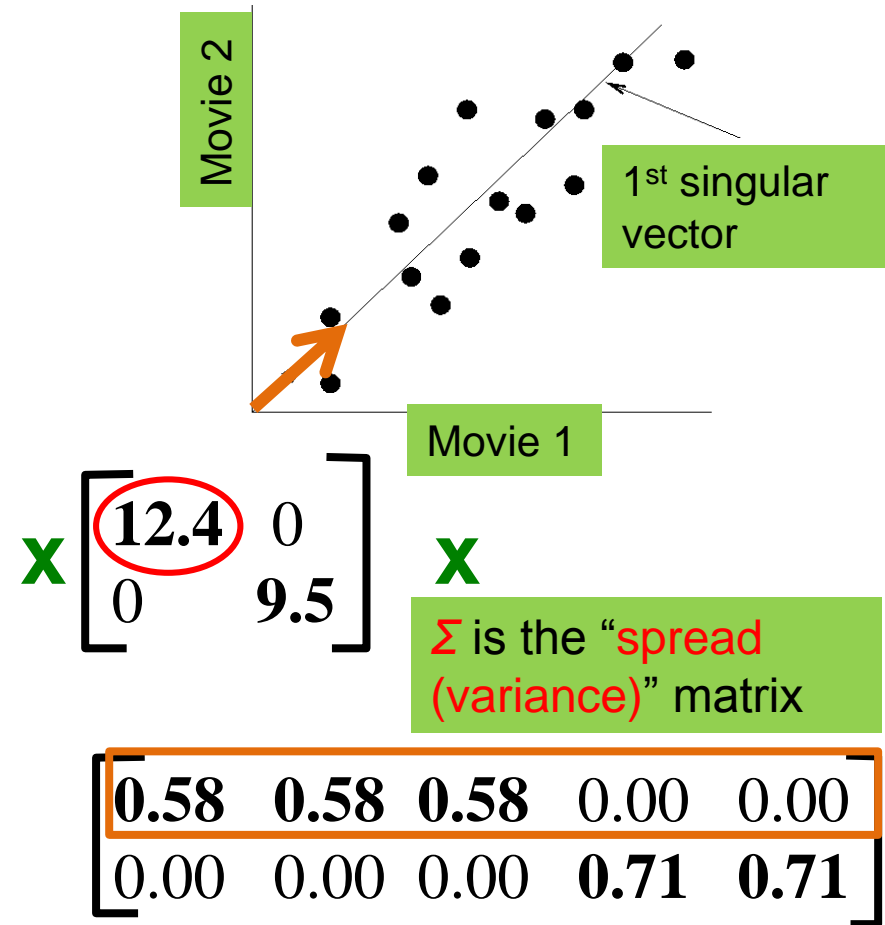
$$= \begin{bmatrix} 0.14 & 0.00 \\ 0.42 & 0.00 \\ 0.56 & 0.00 \\ 0.70 & 0.00 \\ 0.00 & 0.60 \\ 0.00 & 0.75 \\ 0.00 & 0.30 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.71 & 0.71 \end{bmatrix}$$

$\Sigma$  is the “concept strength” matrix

# SVD – Example: Users-to-Movies

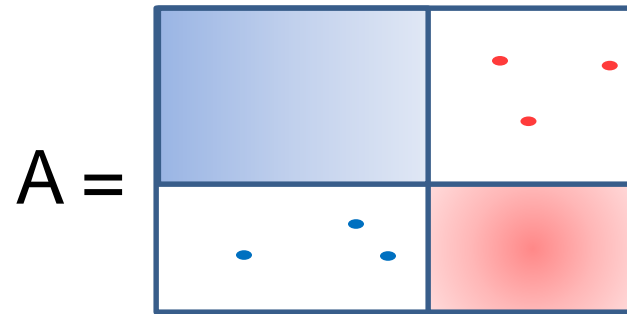
- $A = U \Sigma V^T$  - example: Users to Movies

	Matrix	Alien	Serenity	Casablanca	Amelie	
	1	1	1	0	0	
	3	3	3	0	0	
SciFi	4	4	4	0	0	
	5	5	5	0	0	=
	0	0	0	4	4	$\begin{bmatrix} 0.14 & 0.00 \\ 0.42 & 0.00 \\ 0.56 & 0.00 \\ 0.70 & 0.00 \\ 0.00 & 0.60 \\ 0.00 & 0.75 \\ 0.00 & 0.30 \end{bmatrix}$
Romance	0	0	0	2	2	



# An (more realistic) example

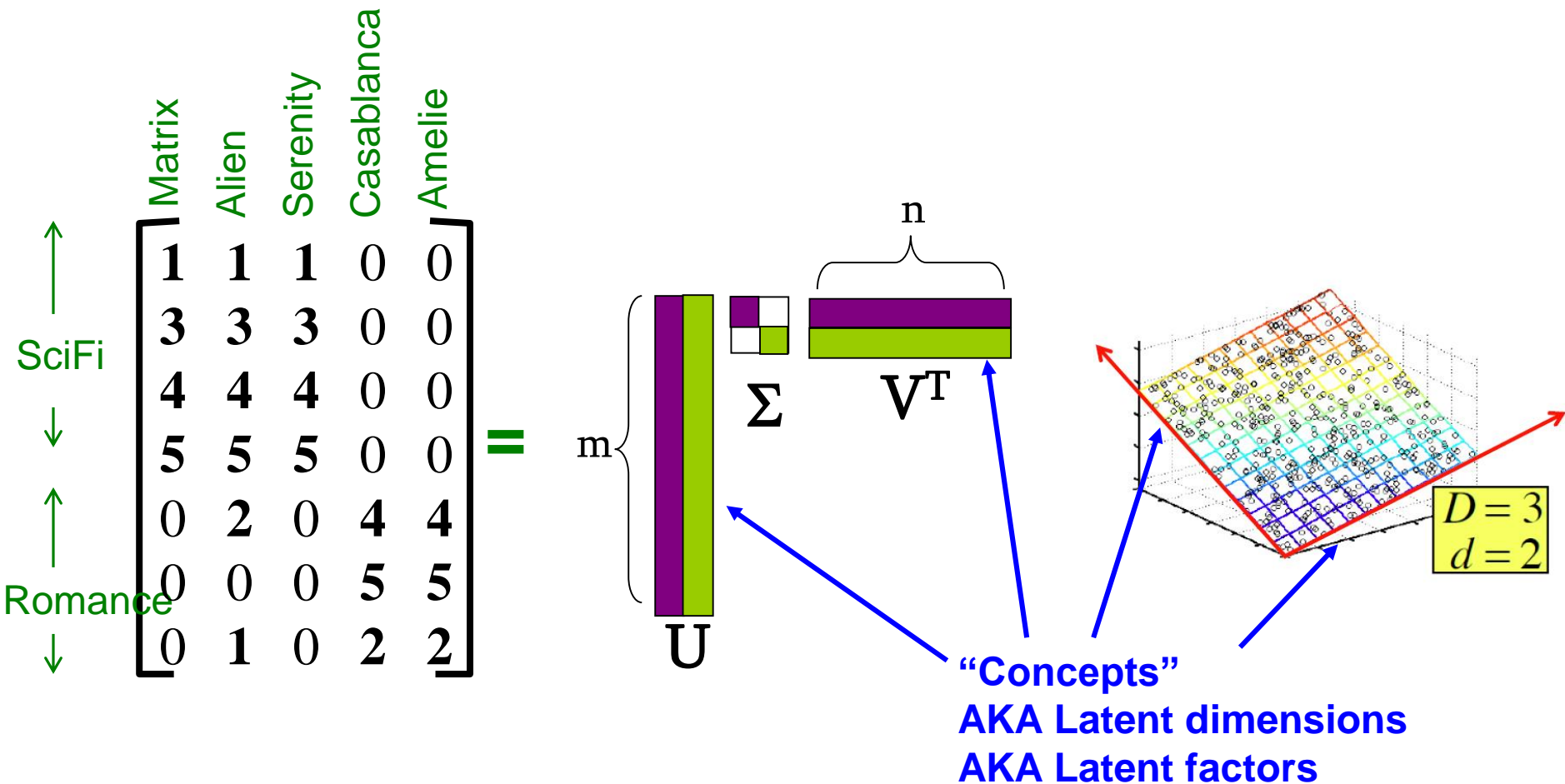
- User-Movie matrix



- There are two prototype users and movies but they are **noisy**

# 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

	Matrix	Alien	Serenity	Casablanca	Amelie										
↑	<b>1</b>	<b>1</b>	<b>1</b>	0	0	<b>=</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>			
SciFi	<b>3</b>	<b>3</b>	<b>3</b>	0	0								<b>0.13</b>	-0.02	-0.01
↓	<b>4</b>	<b>4</b>	<b>4</b>	0	0								<b>0.41</b>	-0.07	-0.03
↑	<b>5</b>	<b>5</b>	<b>5</b>	0	0								<b>0.55</b>	-0.09	-0.04
Romance	0	<b>2</b>	0	<b>4</b>	<b>4</b>								<b>0.68</b>	-0.11	-0.05
↓	0	0	0	<b>5</b>	<b>5</b>								<b>0.15</b>	<b>0.59</b>	<b>0.65</b>
	0	<b>1</b>	0	<b>2</b>	<b>2</b>	<b>0.07</b>	<b>0.73</b>	<b>-0.67</b>							
	0	<b>1</b>	0	<b>2</b>	<b>2</b>	<b>0.07</b>	<b>0.29</b>	<b>0.32</b>							
							<b>X</b>			<b>X</b>					
							<b>12.4</b>	0	0						
							0	<b>9.5</b>	0						
							0	0	<b>1.3</b>						
							<b>X</b>								
							<b>0.56</b>	<b>0.59</b>	<b>0.56</b>	0.09	0.09				
							-0.12	0.02	-0.12	<b>0.69</b>	<b>0.69</b>				
							0.40	<b>-0.80</b>	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-concept	Romance-concept	
↑ SciFi ↓	1	1	1	0	0	=	0.13	-0.02	-0.01
	3	3	3	0	0		0.41	-0.07	-0.03
	4	4	4	0	0		0.55	-0.09	-0.04
↓ Rom ↑	5	5	5	0	0		0.68	-0.11	-0.05
	0	2	0	4	4		0.15	0.59	0.65
	0	0	0	5	5		0.07	0.73	-0.67
	0	1	0	2	2	0.07	0.29	0.32	

The first two vectors are more or less unchanged

x	<table border="0"> <tr> <td style="padding: 5px;">12.4</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;">9.5</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;">1.3</td> </tr> </table>	12.4	0	0	0	9.5	0	0	0	1.3	x
12.4	0	0									
0	9.5	0									
0	0	1.3									

<table border="0"> <tr> <td style="padding: 5px;">0.56</td> <td style="padding: 5px;">0.59</td> <td style="padding: 5px;">0.56</td> <td style="padding: 5px;">0.09</td> <td style="padding: 5px;">0.09</td> </tr> <tr> <td style="padding: 5px;">-0.12</td> <td style="padding: 5px;">0.02</td> <td style="padding: 5px;">-0.12</td> <td style="padding: 5px;">0.69</td> <td style="padding: 5px;">0.69</td> </tr> <tr> <td style="padding: 5px;">0.40</td> <td style="padding: 5px;">-0.80</td> <td style="padding: 5px;">0.40</td> <td style="padding: 5px;">0.09</td> <td style="padding: 5px;">0.09</td> </tr> </table>	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
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

The third vector has a very low singular value

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



# SVD - Interpretation #1

‘**movies**’, ‘**users**’ and ‘**concepts**’:

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

# Rank-k approximation

- In this User-Movie matrix

$$A = \begin{array}{|c|c|} \hline \text{blue block} & \text{white block with 3 red dots} \\ \hline \text{white block with 3 blue dots} & \text{red block} \\ \hline \end{array}$$

- We have more than two singular vectors, but the **strongest** ones are still about the two types.
  - The third models the **noise** in the data
- By keeping the two **strongest singular vectors** we obtain most of the information in the data.
  - This is a **rank-2 approximation** of the matrix A

# Rank- $k$ approximations ( $A_k$ )

$$\begin{pmatrix} A_k \\ n \times d \end{pmatrix} = \begin{pmatrix} U_k \\ n \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \\ k \times k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ k \times d \end{pmatrix}$$

$U_k$  ( $V_k$ ): orthogonal matrix containing the top  $k$  left (right) singular vectors of  $A$ .

$\Sigma_k$ : diagonal matrix containing the top  $k$  singular values of  $A$

$A_k$  is an approximation of  $A$

$A_k$  is the **best** approximation of  $A$

# SVD as an optimization

- The **rank-k approximation** matrix  $A_k$  produced by the top-k singular vectors of A **minimizes** the **Frobenious norm** of the difference with the matrix A

$$A_k = \arg \max_{B: \text{rank}(B)=k} \|A - B\|_F^2$$

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

**Explanation:** The  $(i, j)$  cell in  $A_k$  is close on average with the  $(i, j)$  cell of A

# What does this mean?

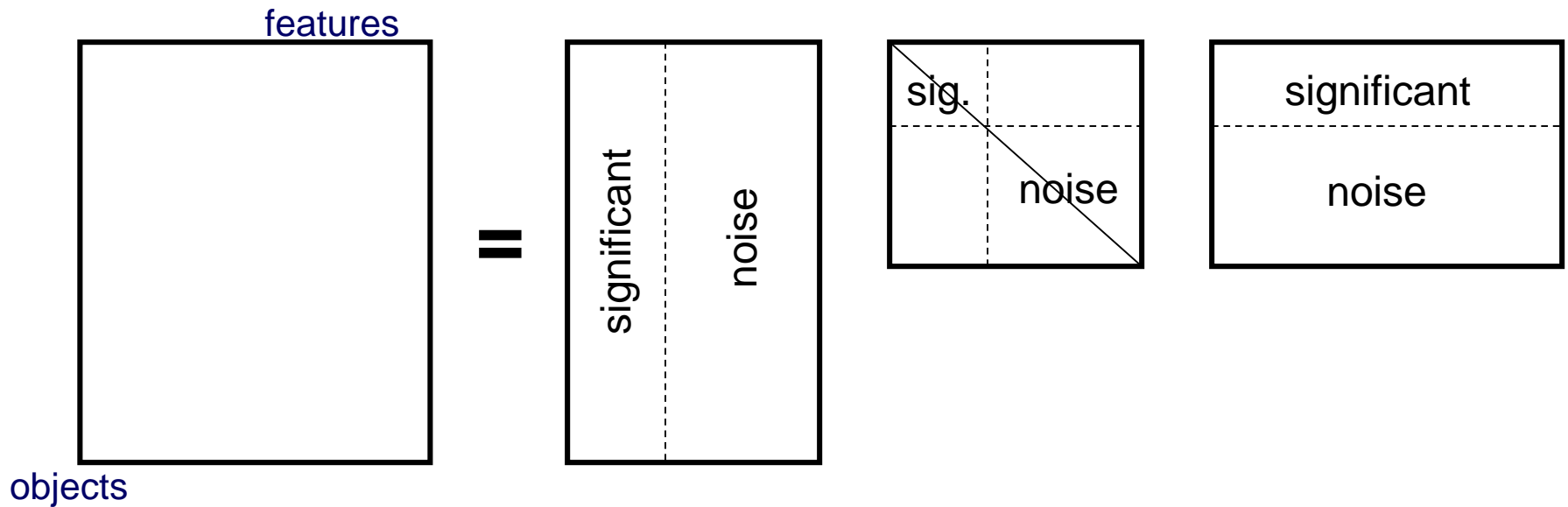
- We can **project** the row (and column) vectors of the matrix  $A$  into a **k-dimensional space** and preserve most of the information
- (**Ideally**) The  $k$  dimensions reveal **latent features/aspects/topics** of the term (document) space.
- (**Ideally**) The  $A_k$  approximation of matrix  $A$ , contains all the **useful information**, and what is discarded is noise

# Latent factor model

- Rows (columns) are linear combinations of  $k$  latent factors
  - E.g., in our extreme document example there are two factors
- Some noise is added to this rank- $k$  matrix resulting in higher rank
- SVD retrieves the latent factors (hopefully).

# SVD and Rank-**k** approximations

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$



# Example

## More details

- **Q:** How exactly is dim. reduction done?
- **A:** Compute SVD

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



# Example

## More details

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

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

# Example

## More details

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

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

# Example

## More details

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

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

# Example

## More details

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

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

# Example

## 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"

# Application: Recommender systems

- Data: Users rating movies
  - Sparse and often noisy
- Assumption: There are  $k$  basic **user profiles**, and each user is a **linear combination** of these profiles
  - E.g., action, comedy, drama, romance
  - Each user is a weighted combination of these profiles
  - The “**true**” matrix has rank  $k$
- If we had the matrix  $A$  with all ratings of all users for all movies, the matrix  $A_k$  would tell us the **true preferences** of the users for the movies

# Model-based Recommendation Systems

- What we observe is a **noisy**, and **incomplete** version of this matrix  $\tilde{A}$
- Given matrix  $\tilde{A}$  and we would like to get the missing ratings that  $A_k$  would produce
- **Algorithm**: compute the **rank-k approximation**  $\tilde{A}_k$  of and matrix  $\tilde{A}$  predict for user  $u$  and movie  $m$ , the value  $\tilde{A}_k[m, u]$ .
  - The rank-k approximation  $\tilde{A}_k$  is **provably** close to  $A_k$
- **Model-based** collaborative filtering

# Example

## Missing ratings and noise

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 3 & 3 & 0 & 0 \\ 4 & 4 & 0 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.14 & -0.06 & -0.04 \\ 0.30 & -0.11 & -0.61 \\ 0.43 & -0.16 & 0.76 \\ 0.74 & -0.31 & -0.18 \\ 0.15 & 0.53 & 0.02 \\ 0.07 & 0.70 & -0.03 \\ 0.07 & 0.27 & 0.01 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.51 & 0.66 & 0.44 & 0.23 & 0.23 \\ -0.24 & -0.13 & -0.21 & 0.66 & 0.66 \\ 0.59 & 0.08 & -0.80 & 0.01 & 0.01 \end{bmatrix}$$



# Example

## Missing ratings and noise

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 3 & 3 & 0 & 0 \\ 4 & 4 & 0 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.14 & -0.06 & -0.04 \\ 0.30 & -0.11 & -0.61 \\ 0.43 & -0.16 & 0.76 \\ 0.74 & -0.31 & -0.18 \\ 0.15 & 0.53 & 0.02 \\ 0.07 & 0.70 & -0.03 \\ 0.07 & 0.27 & 0.01 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.51 & 0.66 & 0.44 & 0.23 & 0.23 \\ -0.24 & -0.13 & -0.21 & 0.66 & 0.66 \\ 0.59 & 0.08 & -0.80 & 0.01 & 0.01 \end{bmatrix}$$

The matrix on the left contains missing ratings (0) and noise (2, 1). The matrix in the middle has its last column crossed out with a red X. The matrix on the right has its last element crossed out with a red X. The multiplication symbols are green.

# Example

- Reconstruction of missing ratings

<b>0.96</b>	<b>1.14</b>	<b>0.82</b>	<b>-0.01</b>	<b>-0.01</b>
<b>1.94</b>	<b>2.32</b>	<b>1.66</b>	<b>0.07</b>	<b>0.07</b>
<b>2.77</b>	<b>3.32</b>	<b>2.37</b>	<b>0.08</b>	<b>0.08</b>
<b>4.84</b>	<b>5.74</b>	<b>4.14</b>	<b>-0.08</b>	<b>0.08</b>
<b>0.40</b>	<b>1.42</b>	<b>0.33</b>	<b>4.06</b>	<b>4.06</b>
<b>-0.42</b>	<b>0.63</b>	<b>-0.38</b>	<b>4.92</b>	<b>4.92</b>
<b>0.20</b>	<b>0.71</b>	<b>0.16</b>	<b>2.03</b>	<b>2.03</b>

# SVD and PCA

- PCA is a special case of SVD on the **centered matrix**.
  - A matrix where we subtract the mean from the row.

# Covariance matrix

- Goal: reduce the dimensionality while preserving the “information in the data”
- Information in the data: **variability** in the data
  - We measure variability using the **covariance matrix**.
  - Sample covariance of variables X and Y

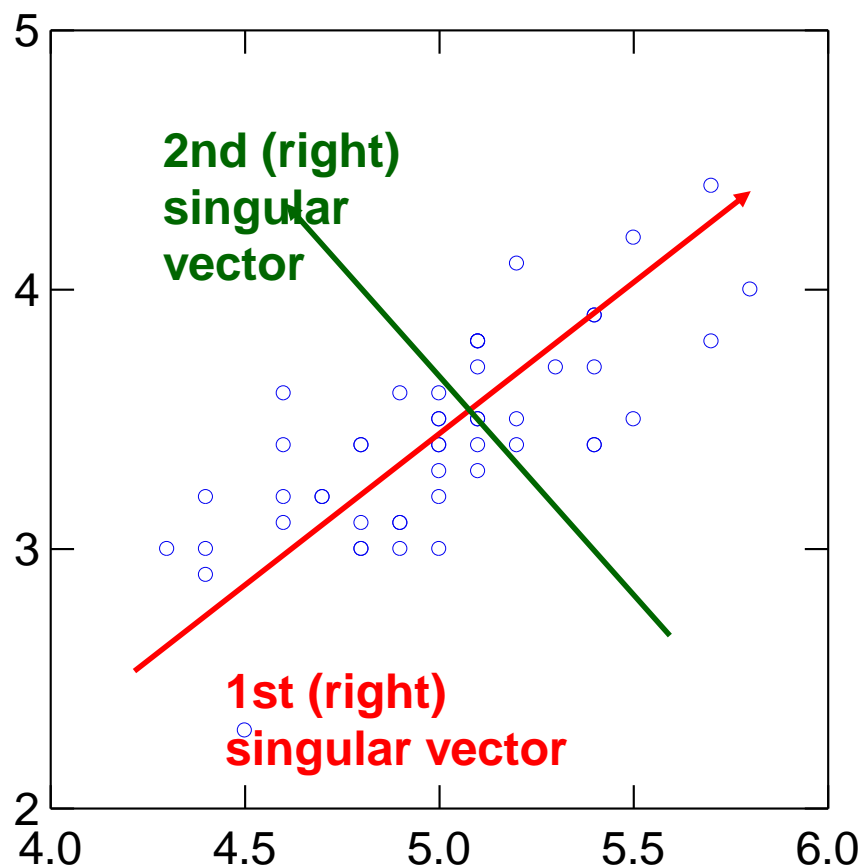
$$\sum_i (x_i - \mu_X)^T (y_i - \mu_Y)$$

- Given matrix **A**, remove the **mean** of each column from the column vectors to get the **centered** matrix **C**
- The matrix  $V = C^T C$  is the **covariance matrix** of the **row** vectors of **A**.

# PCA: Principal Component Analysis

- We will project the rows of matrix **C** into a new set of **attributes** (dimensions) such that:
  - The attributes have **zero covariance** to each other (they are **orthogonal**)
  - Each attribute captures **the most remaining variance** in the data, while orthogonal to the existing attributes
    - The first attribute should capture the most variance in the data
- For matrix **C**, the variance of the rows of **C** when projected to vector  $x$  is given by  $\sigma^2 = ||Cx||^2$ 
  - The **right singular vector of C** maximizes  $\sigma^2$ !

# PCA



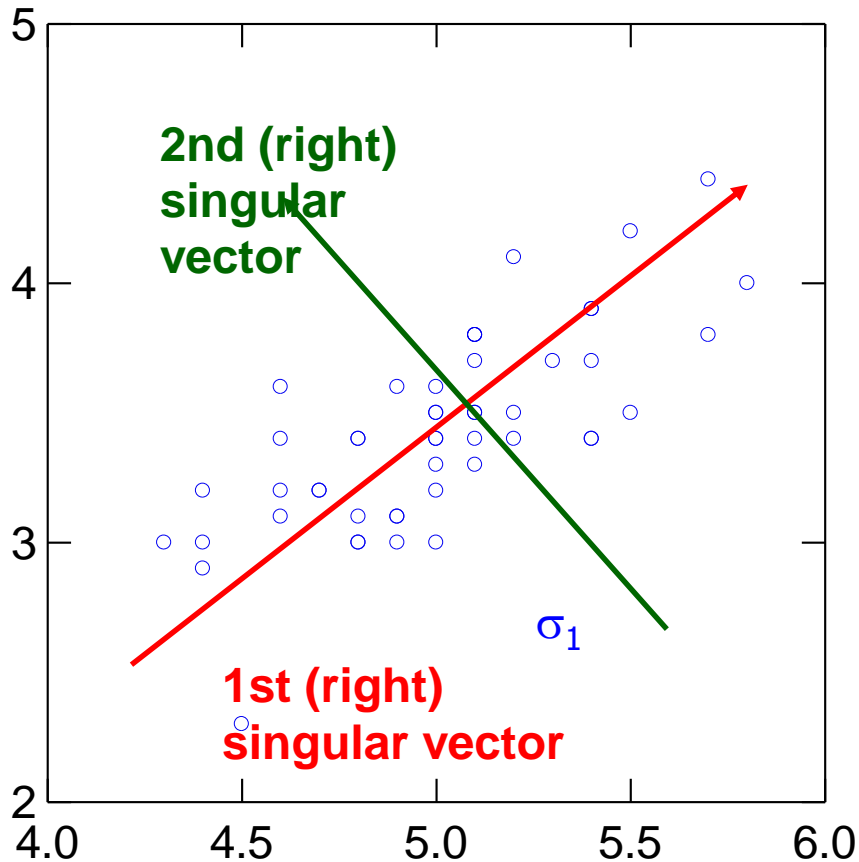
**Input:** 2-d dimensional points

**Output:**

**1st (right) singular vector:**  
direction of maximal variance,

**2nd (right) singular vector:**  
direction of maximal variance,  
after removing the projection of  
the data along the first singular  
vector.

# Singular values



$\sigma_1$ : measures how much of the data variance is explained by the first singular vector.

$\sigma_2$ : measures how much of the data variance is explained by the second singular vector.

# Singular values tell us something about the variance

- The variance in the direction of the **k**-th principal component is given by the corresponding singular value  $\sigma_k^2$
- Singular values can be used to estimate how many components to keep
- **Rule of thumb:** keep enough to explain **85%** of the variation:

$$\frac{\sum_{j=1}^k \sigma_j^2}{\sum_{j=1}^n \sigma_j^2} \approx 0.85$$



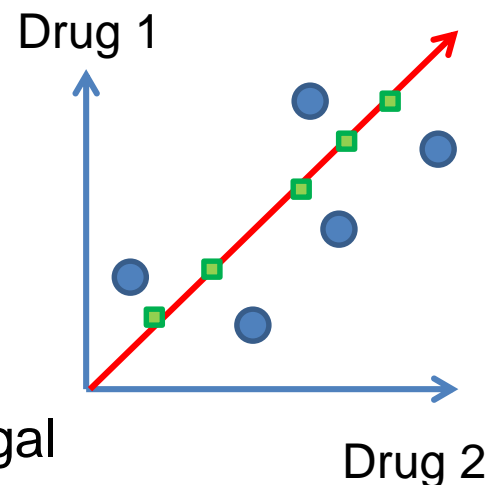
# Example

$$A = \begin{matrix} & \text{drugs} & & \\ & \begin{matrix} \cdots & \cdots & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \cdots & \cdots \end{matrix} & & \\ \begin{matrix} a_{11} \\ \vdots \\ a_{n1} \end{matrix} & \begin{matrix} \cdots \\ \vdots \\ \cdots \end{matrix} & \begin{matrix} a_{1n} \\ \vdots \\ a_{nn} \end{matrix} & \text{students} \\ & \text{legal} & \text{illegal} & \end{matrix}$$

$a_{ij}$ : usage of student  $i$  of drug  $j$

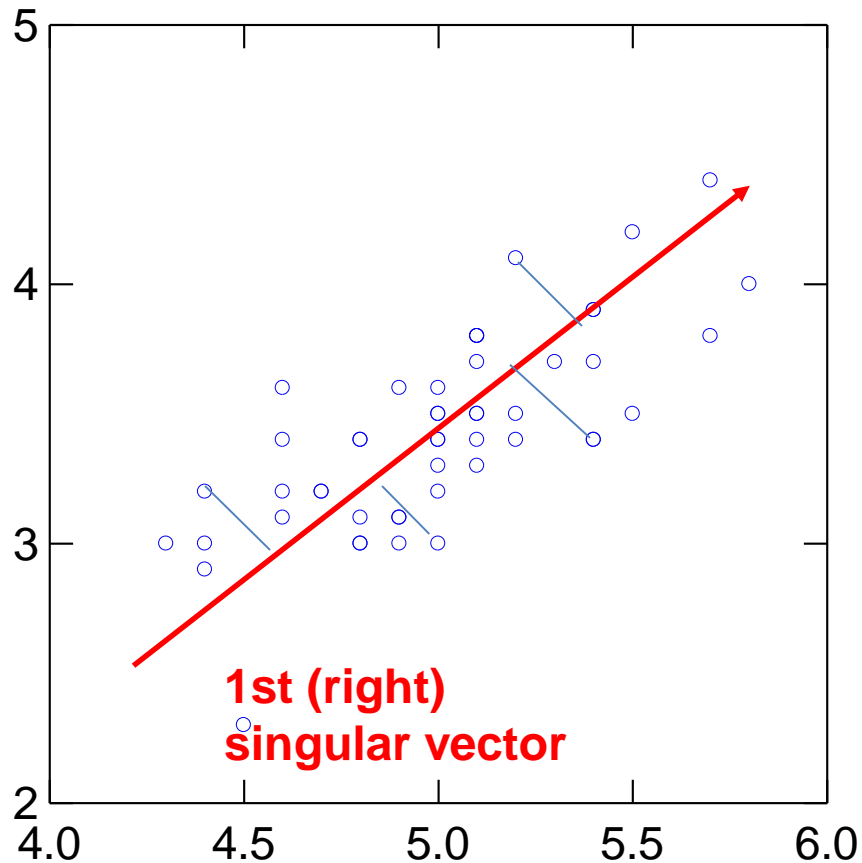
$$A = U\Sigma V^T$$

- First right singular vector  $v_1$ 
  - More or less same weight to all drugs
  - Discriminates heavy from light users
- Second right singular vector
  - Positive values for legal drugs, negative for illegal



# Another property of PCA/SVD

- The chosen vectors are such that minimize the **sum of square differences** between the data vectors and the low-dimensional projections



# Application

- **Latent Semantic Indexing (LSI):**
  - Apply PCA on the **document-term** matrix, and index the k-dimensional vectors
  - When a query comes, **project** it onto the k-dimensional space and compute **cosine similarity** in this space
  - Principal components capture main **topics**, and enrich the document representation

**SVD is “the Rolls-Royce and the Swiss Army Knife of Numerical Linear Algebra.”\***

**\*Dianne O’Leary, MMDS ’06**

# Computation of eigenvectors

- Consider a symmetric square matrix  $M$
- Power-method:
  - Start with the vector  $v$  of all 1's
  - Compute  $v = Mv$
  - Normalize by the length of  $v$
  - Repeat until the vector does not change
- This will give us the first eigenvector.
- The first eigenvalue is  $\lambda = v^T Mv$
- For the second one, compute the first eigenvector of the matrix  $M^* = M - \lambda v v^T$

# Singular Values and Eigenvalues

- The **left singular vectors** of  $A$  are also the eigenvectors of  $AA^T$
- The **right singular vectors** of  $A$  are also the eigenvectors of  $A^T A$
- The **singular values** of matrix  $A$  are also the square roots of eigenvalues of  $AA^T$  and  $A^T A$

# Computing singular vectors

- Compute the eigenvectors and eigenvalues of the matrices  $MM^T$  and  $M^T M$