# Mining Massive Datasets Dimensionality Reduction SVD&CUR
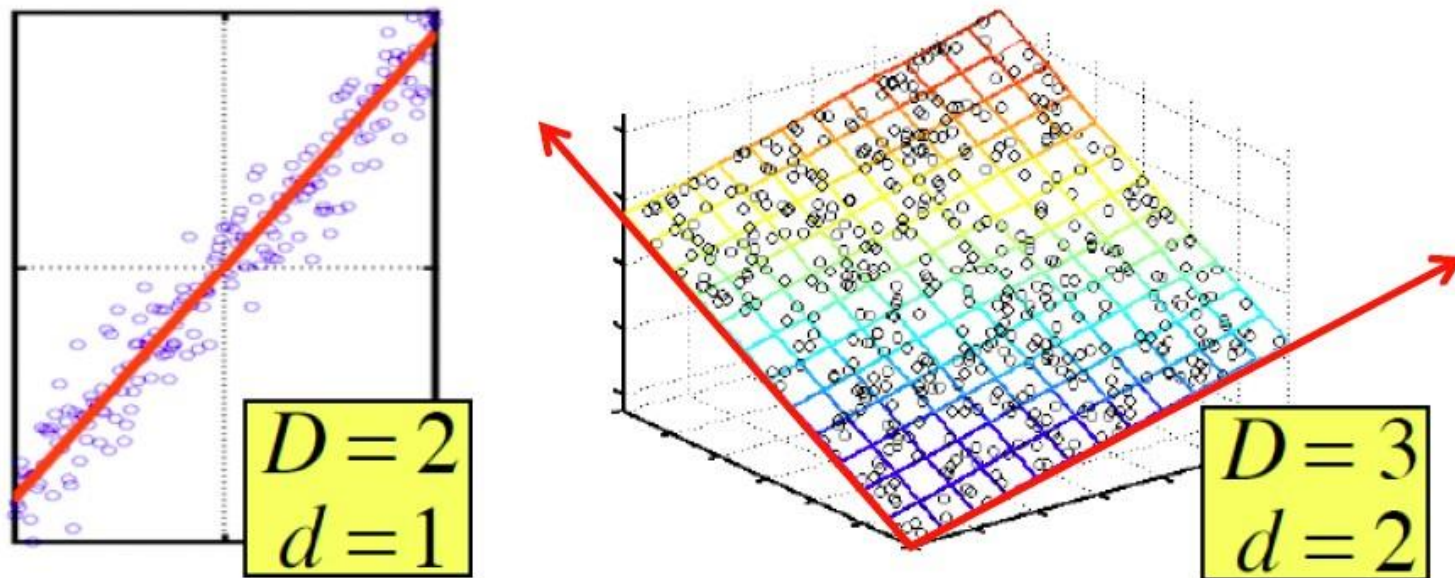
## ■ Compress / reduce dimensionality:

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

| day customer | 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]

$D = 2$
$d = 1$

$D = 3$
$d = 2$

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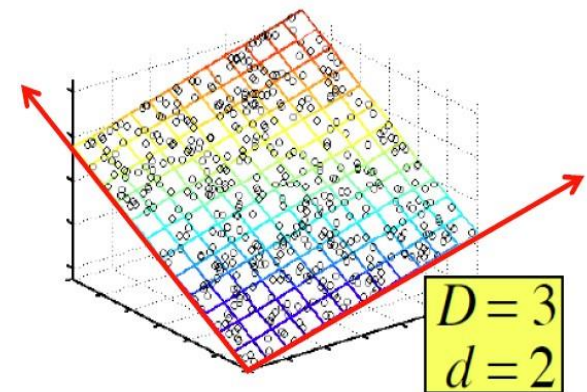- Axes of this subspace are effective representation of the data

- 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?

- **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**
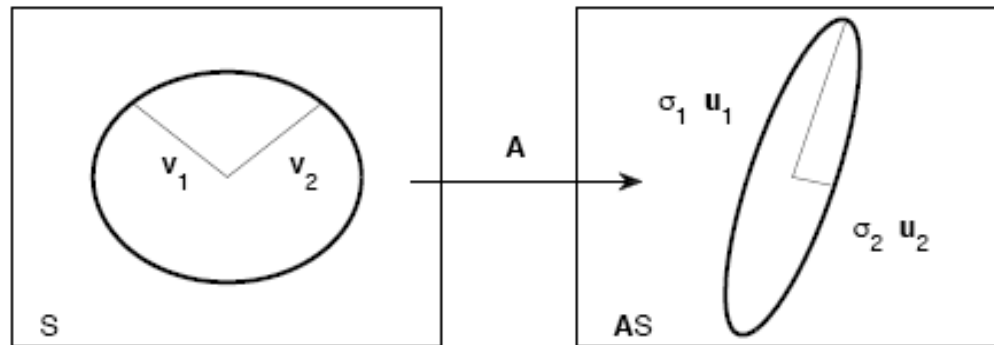
$D = 3$
$d = 2$

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

- **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)

• 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

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

Given a n x n matrix A$^{\mathsf{T}}$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:

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

where $A^\top A$ is the covariance matrix of the data

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

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

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

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

# SVD-decomposition

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

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

where $A^T A$ is the covariance matrix of the data
Objective: maximize variance subject to constraint $v^T v = 1$.
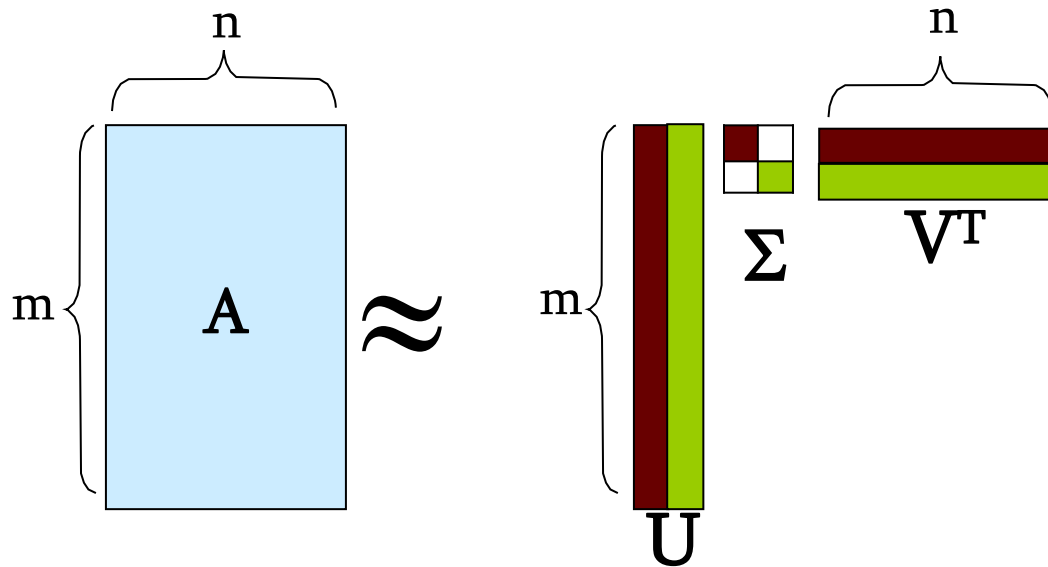
Maximize $\quad f = v^T A^T A v - \sigma(v^T v - 1)$

σ is the Lagrange multiplier, Differentiating with respect to v yields Eigenvalue equation:

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

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

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



$$\sigma_i \quad \dots \quad \text{scalar}$$
$$\mathbf{u}_i \quad \dots \quad \text{vector}$$
$$\mathbf{v}_i \quad \dots \quad \text{vector}$$
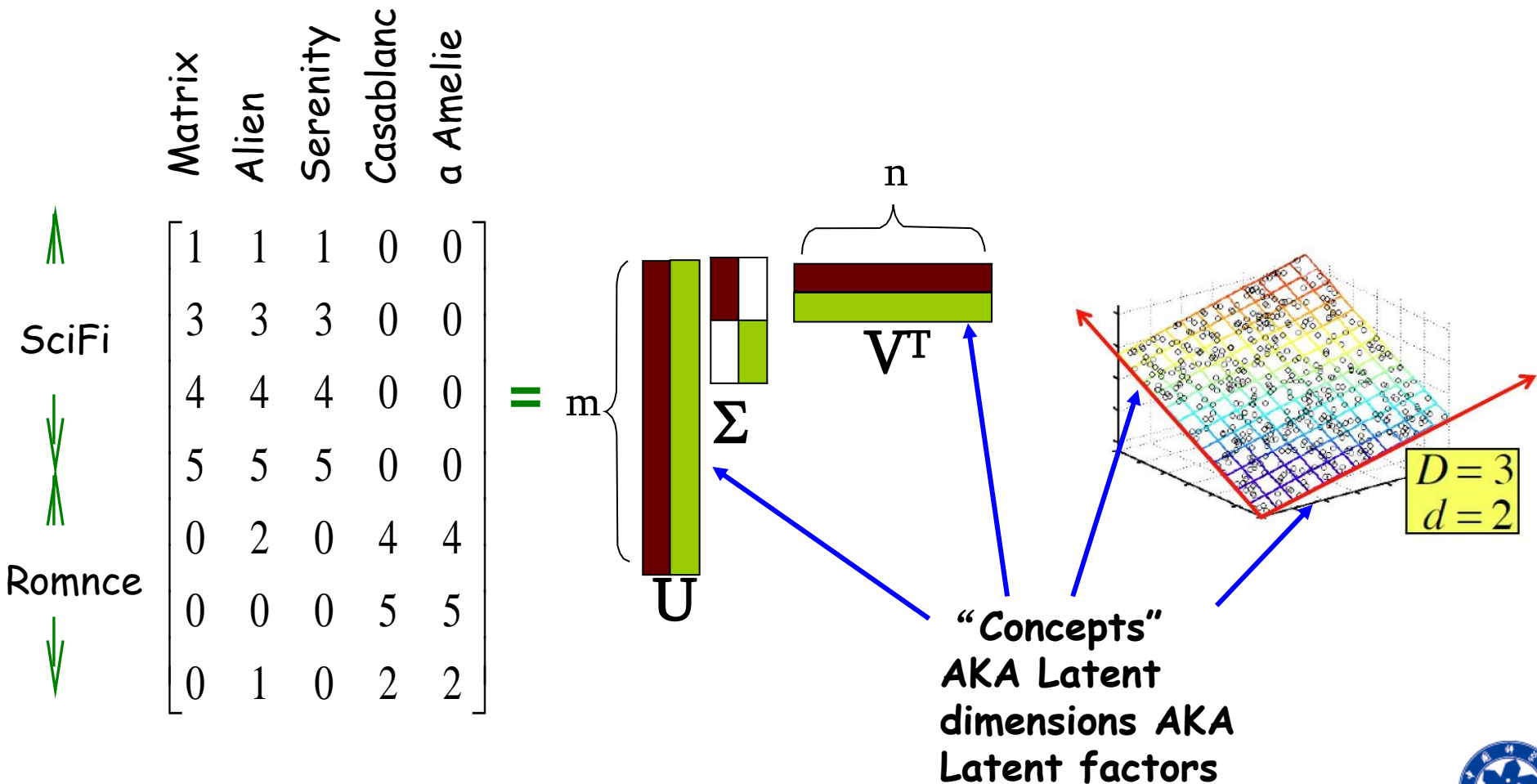
It is **always** possible to decompose a real matrix $A$ into $A = U\sum V^T$ , where

- $U$, $\Sigma$, $V$: 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 ($\sigma_1 \geq \sigma_2 \geq \ldots \geq 0$)

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



$$\begin{array}{c} \text{Matrix} \quad \text{Alien} \quad \text{Serenity} \quad \text{Casablanca} \quad \text{Amelie} \end{array}$$

SciFi

Romnce

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

= m

U   Σ   $V^T$   n

$D = 3$
$d = 2$

"Concepts"
AKA Latent
dimensions AKA
Latent factors

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

SciFi

Romnce

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

Matrix, Alien, Serenity, Casablanc, Amelie

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

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

SciFi-concept

Romance-concept

$$
\begin{array}{c}
\text{SciFi} \\
\\
\text{Romnce}
\end{array}
\quad
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
3 & 3 & 3 & 0 & 0 \\
4 & 4 & 4 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 2 & 0 & 4 & 4 \\
0 & 0 & 0 & 5 & 5 \\
0 & 1 & 0 & 2 & 2
\end{bmatrix}
=
\begin{bmatrix}
0.13 & 0.02 & -0.01 \\
0.41 & 0.07 & -0.03 \\
0.55 & 0.09 & -0.04 \\
0.68 & 0.11 & -0.05 \\
0.15 & -0.59 & 0.65 \\
0.07 & -0.73 & -0.67 \\
0.07 & -0.29 & 0.32
\end{bmatrix}
\times
\begin{bmatrix}
12.4 & 0 & 0 \\
0 & 9.5 & 0 \\
0 & 0 & 1.3
\end{bmatrix}
\times
$$

(columns: Matrix, Alien, Serenity, Casablanc, Amelie)

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

- **A = UΣVᵀ - example: U is "user-to-concept" similarity matrix**

SciFi-concept    Romance-concept

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

Movies: Matrix, Alien, Serenity, Casablanc, a Amelie

SciFi

Romnce

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

SciFi-concept

"strength" of the SciFi-concept

|  | Matrix | Alien | Serenity | Casablanc | a Amelie |
|---|---|---|---|---|---|
| SciFi | 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 |
| Romnce | 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** |

**x**

| | | |
|---|---|---|
| **12.4** | 0 | 0 |
| 0 | **9.5** | 0 |
| 0 | 0 | **1.3** |

**x**

| | | | | |
|---|---|---|---|---|
| **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 |

- **A = UΣVᵀ - example: V is "movie-to-concept" similarity matrix**



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

Columns: Matrix, Alien, Serenity, Casablanc, Amelie

SciFi / Romnce

**SciFi-concept**

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

**SciFi-concept**

## '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 gives 'best' axis to project on:**

  - '**best**' = min sum of squares of projection errors
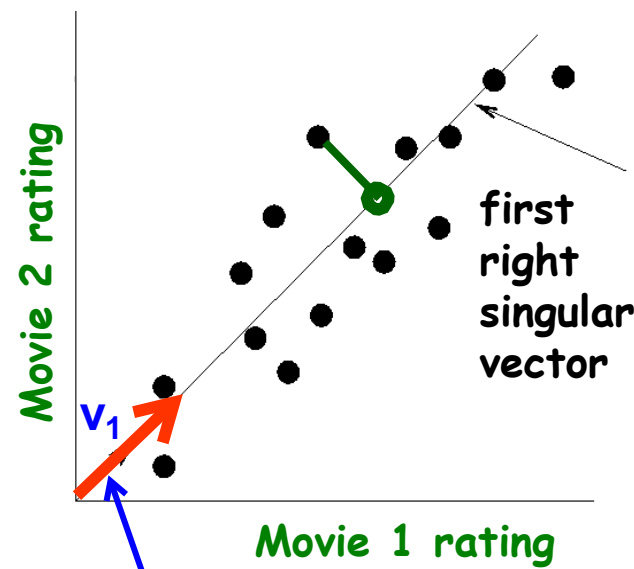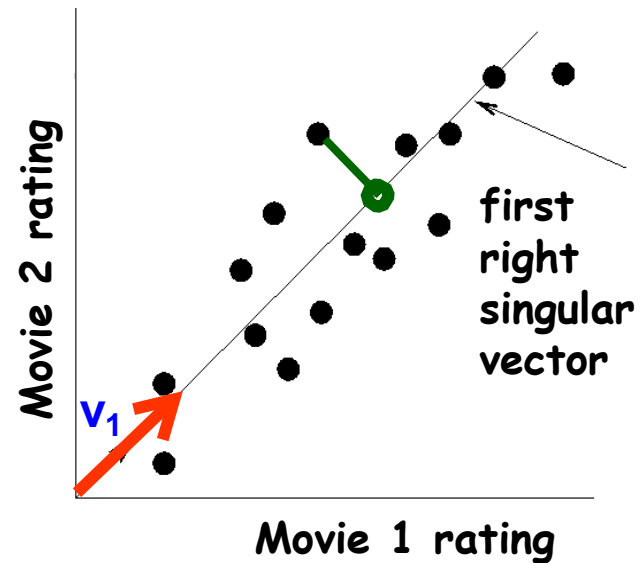
  - ## In other words,
    minimum reconstruction error



Movie 2 rating

Movie 1 rating

first right singular vector

$v_1$

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



first right singular vector

Movie 2 rating

Movie 1 rating

$v_1$

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

# SVD-Interpretation #2

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



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

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

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

- **A = UΣVᵀ - example:**

  - **UΣ:**  Gives the coordinates of the points in the projection axis



Movie 2 rating

first right singular vector

$v_1$

Movie 1 rating

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

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

$$\begin{bmatrix} 1.61 & 0.19 & -0.01 \\ 5.08 & 0.66 & -0.03 \\ 6.82 & 0.85 & -0.05 \\ 8.43 & 1.04 & -0.06 \\ 1.86 & -5.60 & 0.84 \\ 0.86 & -6.93 & -0.87 \\ 0.86 & -2.75 & 0.41 \end{bmatrix}$$

## More details
- Q: How exactly is dim. reduction done?

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

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

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

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

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

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

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

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

# More details

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

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

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

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

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

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

## 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:**

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

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

is "small"

# SVD-Best Low Rank Approx.

A = U Sigma V^T

**B is best approximation of A**

B = U Sig ma V^T

- <u>Theorem:</u> Let A = U $\Sigma$ V$^T$      ($\sigma_1 \geq \sigma_2 \geq ...$, rank($A$)=$r$) then B = U S V$^T$

  - S = diagonal nxn matrix where $s_i = \sigma_i$ (i=1...k) else $s_i = 0$

## is a best rank-k approximation to $A$:

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

$$\begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & \\ \vdots & \vdots & \ddots & \\ x_{m1} & & & x_{mn} \end{pmatrix}_{m \times n} = \overset{U}{\begin{pmatrix} u_{11} & \cdots & \\ \vdots & \ddots & \\ u_{m1} & & \end{pmatrix}}_{m \times r} \overset{\Sigma}{\begin{pmatrix} \sigma_1 & 0 & \cdots \\ 0 & \ddots & \\ \vdots & & \end{pmatrix}}_{r \times r} \overset{V^T}{\begin{pmatrix} v_{11} & \cdots & v_{1n} \\ \vdots & \ddots & \\ & & \end{pmatrix}}_{r \times n}$$

## - We will need 2 facts:

  - $\left\| M \right\|_F = \sum_i (q_{ii})^2$ where $\boldsymbol{M} = \boldsymbol{P}\,\boldsymbol{Q}\,\boldsymbol{R}$ is SVD of $\boldsymbol{M}$
  - $\boldsymbol{U}\,\Sigma V^T - \boldsymbol{U}\,\boldsymbol{S}\,V^T = \boldsymbol{U}\,(\Sigma - \boldsymbol{S})\,V^T$

- ## **We will need 2 facts:**
  - $\left\|M\right\|_F = \sum_k (q_{kk})^2$ 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

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

**We apply:**
-- P column orthonormal
-- R row orthonormal
-- Q is diagonal

- ■ **A = U $\Sigma$ V$^T$ , B = U S V$^T$** ($\sigma_1 \geq \sigma_2 \geq \ldots \geq 0$, rank($A$)=$r$)
  - ▪ **S** = diagonal $n \times n$ matrix where $s_i = \sigma_i$ ($i=1 \ldots k$) else $s_i=0$

**then $B$** is solution to $\min_B \|A\text{-}B\|_F$, rank($B$)=$k$

- ■ **Why?**

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

**We used: U $\Sigma$V$^T$ - U S V$^T$ = U ($\Sigma$- S) V$^T$**

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

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

**Equivalent:**
**'spectral decomposition' of the matrix:**

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} | & | \\ | & | \\ U_1 & U_2 \\ | & | \\ | & | \end{bmatrix} \times \begin{bmatrix} \sigma_1 & \oslash \\ \oslash & \sigma_2 \end{bmatrix} \times \begin{bmatrix} \underline{\qquad} V_1 \underline{\qquad} \\ \underline{\qquad} V_2 \underline{\qquad} \end{bmatrix}$$

# Equivalent:
## 'spectral decomposition' of the matrix:

$$\underset{n}{\overset{\longleftarrow m \longrightarrow}{\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix}}}$$

$\longleftarrow$ k terms $\longrightarrow$

$= \quad \sigma_1 \quad u_1 \quad v^T_1 \quad + \quad \sigma_2 \quad u_2 \quad v^T_2 \quad + ...$

$n \times 1 \qquad 1 \times m$

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

**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.

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

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

$$
\begin{array}{c} \xleftarrow{\hspace{1cm}} m \xrightarrow{\hspace{1cm}} \end{array}
$$

$$
n \left|
\begin{matrix}
1 & 1 & 1 & 0 & 0 \\
3 & 3 & 3 & 0 & 0 \\
4 & 4 & 4 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 2 & 0 & 4 & 4 \\
0 & 0 & 0 & 5 & 5 \\
0 & 1 & 0 & 2 & 2
\end{matrix}
\right|
= \sigma_1 \quad u_1 \quad v^T_1 \quad + \quad \sigma_2 \quad u_2 \quad v^T_2 \quad +\ldots
$$

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

- **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: 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

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

- **SVD gives us:**
  - $A = U \, \Sigma V^T$
- **Eigen-decomposition:**
  - $A = X \, \Lambda X^T$
    - A is symmetric
    - U, V, X are orthonormal ($\mathbf{U^T U = I}$),
    - $\Lambda$, $\Sigma$ are diagonal
- **What is:**
  - $\mathbf{AA^T = U\Sigma V^T (U\Sigma V^T)^T = U\Sigma V^T (V\Sigma U^T) = U\Sigma\Sigma^T \, U^T}$
  - $\mathbf{A^T A = V \, \Sigma \, U^T (U\Sigma V^T) = V \, \Sigma\Sigma^T V^T}$

Shows how to compute SVD using eigenvalue decomposition!

$X \, \Lambda \; X^T$

$X \, \Lambda \; X^T$

So, $\lambda_i = \sigma_i^2$

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

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

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

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

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

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

**Project into concept space:**
Inner product with each 'concept' vector $v_i$

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

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

**Project into concept space:** Inner product with each 'concept' vector $v_i$

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

## E.g.:

$$q = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.12 \\ 0.59 & -0.02 \\ 0.56 & 0.12 \\ 0.09 & -0.69 \\ 0.09 & -0.69 \end{bmatrix} = \begin{bmatrix} 2.8 & 0.6 \end{bmatrix}$$

(columns: Matrix, Alien, Serenity, Casablanca, Amelie)

SciFi-concept

**movie-to-concept similarities (V)**

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

$$\mathbf{d_{concept}} = \mathbf{d\ V}$$

**E.g.:**

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

SciFi-concept

movie-to-concept similarities (V)

- **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**!



$$d= \begin{bmatrix} 0 & 4 & 5 & 0 & 0 \end{bmatrix} \quad\dashrightarrow\quad \begin{bmatrix} 2.8 & 0.6 \end{bmatrix}$$

$$q= \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \end{bmatrix} \quad\dashrightarrow\quad \begin{bmatrix} 5.2 & 0.4 \end{bmatrix}$$

SciFi-concept

**Zero ratings in common**

Similarity ≠ 0

# 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:**

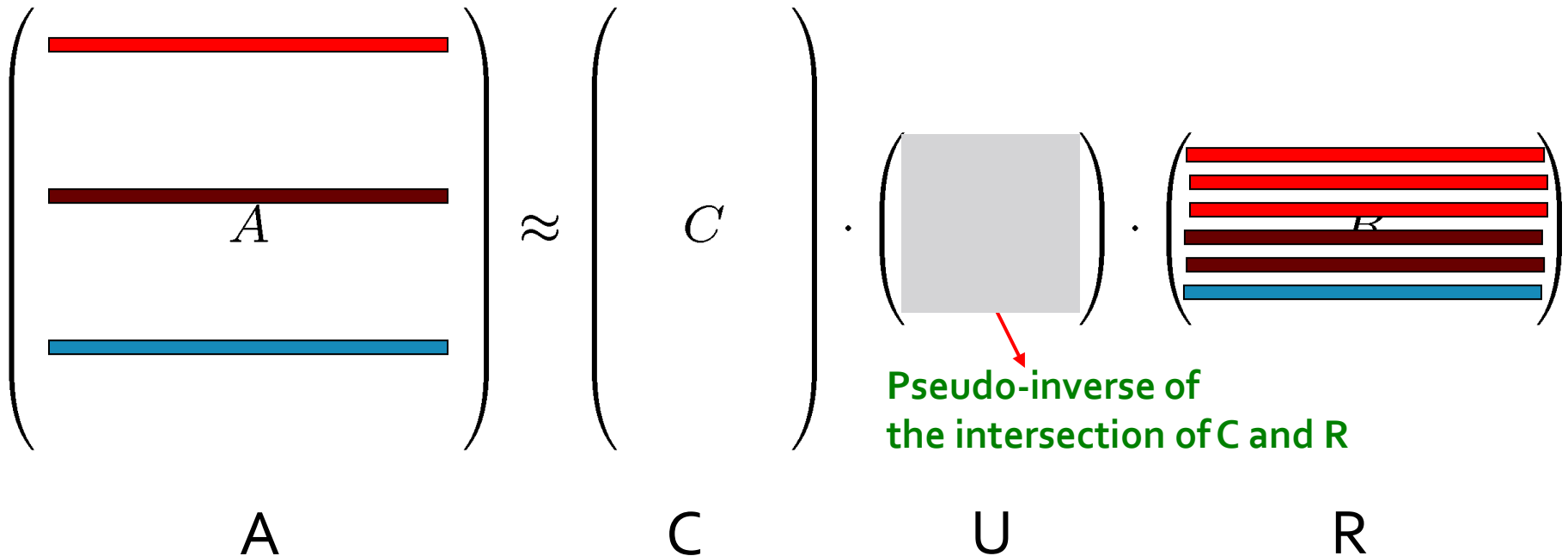  - Singular vectors are **dense!**

$$\text{(matrix)} = U \quad \Sigma \quad V^T$$

# CUR Decomposition

Frobenius norm:

$$\|X\|_F = \Sigma_{ij} \; X_{ij}{}^2$$

- **Goal: Express A as a product of matrices C,U,R Make ‖A−C·U·R‖$_F$ small**



$$A \approx C \cdot U \cdot R$$

A          C          U          R

Frobenius norm:

$$\|X\|_F = \Sigma_{ij} \, X_{ij}{}^2$$

- **Goal: Express A as a product of matrices C,U,R Make ‖A-C·U·R‖$_F$ small**
- **"Constraints" on C and R:**



$$A \approx C \cdot U \cdot R$$

Pseudo-inverse of
the intersection of C and R

A        C        U        R

- **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·n**) time achieves

- $\|A\text{-}CUR\|_F \leq \|A\text{-}A_k\|_F + \varepsilon \|A\|_F$

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

- **O(k log(1/$\delta$ ) /$\varepsilon^2$)** columns, and

- **O($k^2$ log$^3$(1/$\delta$ )/$\varepsilon^6$ )** rows

**In practice:**
Pick 4*k*
cols/rows

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

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

| A | ≈ | W | R |
|---|---|---|---|
|   |   | C |   |

$$U = W^+$$

**Why pseudoinverse works?**
W = X Z Y then W$^{-1}$ = X$^{-1}$ Z$^{-1}$ Y$^{-1}$
Due to orthonomality X$^{-1}$=X$^T$ and Y$^{-1}$=Y$^T$
  Since Z is diagonal Z = $1/Z_{ii}$
**Thus**, if **W** is nonsingular, pseudoinverse is the true inverse

+ 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

Actual column
Singular vector

■ **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

$R_d$

$C_d$

$R_s$

$C_s$    Construct a small U

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

sparse and small

Huge but sparse    Big and dense

CUR:  $A = C U R$

dense but small

Huge but sparse    Big but sparse

- **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?

- **Accuracy:**
  - 1 – relative sum squared errors

- **Space ratio:**
  - #output matrix entries / #input matrix entries

- **CPU time**

Sun, Faloutsos: *Less is More: Compact Matrix Decomposition for Large Sparse Graphs*, SDM '07
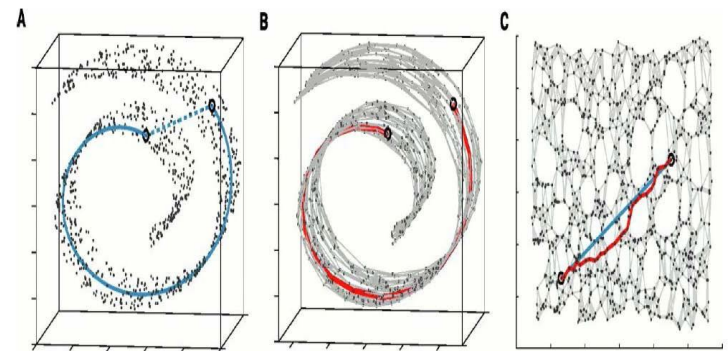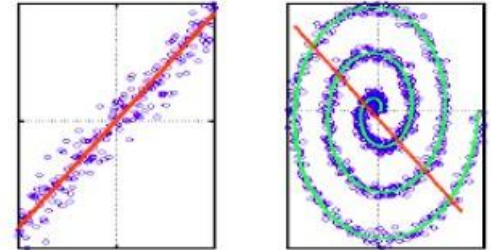
- **SVD is limited to linear projections:**

  

  - 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