

# Machine Learning

## Lecture 10: Dimensionality Reduction & Matrix Factorization

---

Prof. Dr. Stephan Günnemann

Data Analytics and Machine Learning Group  
Technical University of Munich

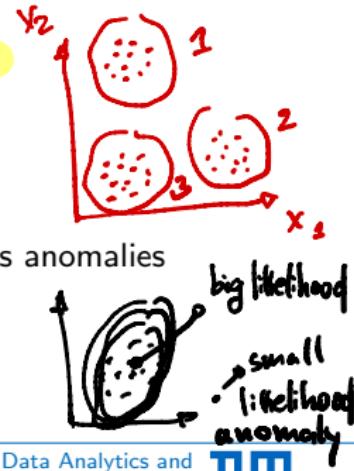
Winter term 2022/2023

## Chapter: Dimensionality Reduction & Matrix Factorization

- 1. Introduction**
2. Principal Component Analysis (PCA)
3. Singular Value Decomposition (SVD)
4. Matrix Factorization
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

# Introduction: Unsupervised Learning

- Supervised learning aims to map inputs to targets with  $y = f(x)$ , or in a probabilistic framework to model  $p(y|x)$
- Unsupervised learning can be seen as **modeling  $p(x)$**
- We are trying to find the (hidden/latent) structure in the data
  - e.g find a latent distribution  $p(z)$  and a generative transformation  $p(x|z)$  to obtain  $p(x) = \int p(x|z)p(z)dz$
  - **Latent  $z$  usually unknown and has to be estimated**
- Examples:
  - Clustering: The cluster label is the latent state
  - Anomaly detection: treat instances with low  $p(x)$  as anomalies



# Introduction: Unsupervised Learning

- Unsupervised learning can be viewed as compression
  - Compress a data point to a single label corresponding to its cluster
  - Compress a data point from a higher to a lower dim. latent space
- Unsupervised learning can be used:
  - As standalone method (e.g. to understand your data, visualization)
  - As pre-processing step (e.g. use cluster label as feature for subsequent classification task; obtain small number of relevant features)
  - To leverage large amounts of unlabeled data for pretraining
- This lecture: Dimensionality Reduction (& Matrix Factorization)

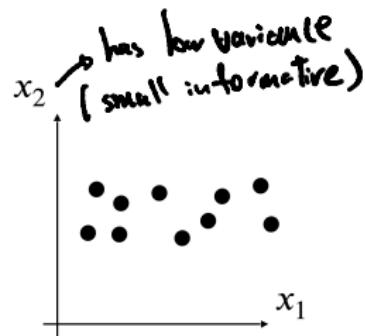
# Dimensionality Reduction: Motivation

- Often data has very many features, i.e. high-dimensional data
- High-dimensional data is challenging:
  - Similarity search/computation is expensive because of high complexity of distance functions
  - Highly correlated dimension could cause trouble for some algorithms
  - Curse of dimensionality: we need exponential amounts of data to characterize the density as the dimensionality goes up
  - It is hard to visualize high-dimensional data
- Often the data lies on a low-dim. manifold, embedded in a high-dim. space
- Goal: Try to reduce the dimensionality while avoiding information loss
- Benefits:
  - Computational or memory savings
  - Uncover the intrinsic dimensionality of the data
  - (more benefits later...)

# Feature (Sub-)Selection

Choose "good" dimensions using a-priori knowledge or appropriate heuristics

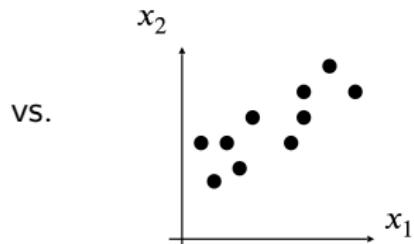
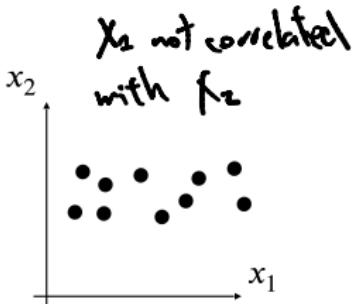
- e.g. remove **low-variance** dimensions
- Depending on the application only a few dimensions might be of interest
  - Example: shoe size interesting for shoe purchases, not so for car purchases
- Advantage:
  - No need for an intensive preprocessing or training phase to determine relevant dimensions
- Disadvantage:
  - Expert knowledge required; misjudgment possible
  - Univariate feature selection ignores correlations



# Beyond Feature (Sub-)Selection

Can we do

- better?
- automatic?

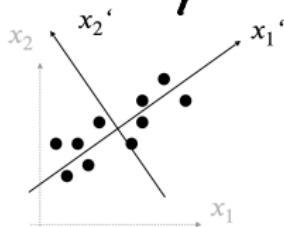
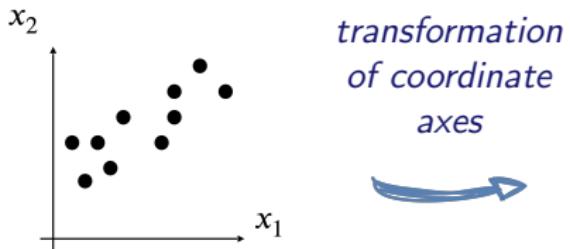


Obviously: Simply discarding whole features not a good idea

- Features are often correlated

now  $x_2'$  is not correlated with  $x_1'$ , so we can discard  $x_2'$

The  $x_2'$  dimension is barely necessary to describe the data in this new coordinate system  
→ can be ignored

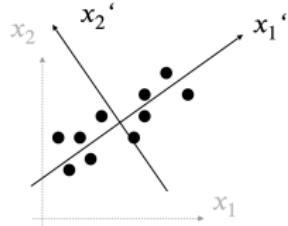
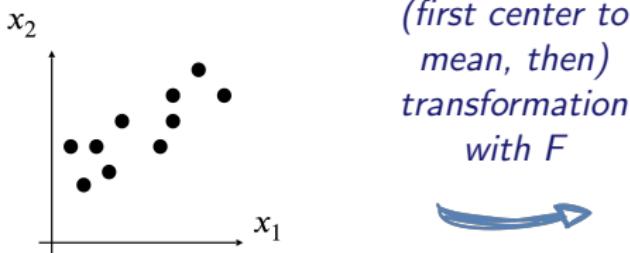


# Dim. Reduction via Linear Transformations

- Use linear transformations to represent data in a different coordinate system
  - **Change of basis** (orthogonal basis transformations) and **potentially discarding dimensions**
$$\begin{matrix} F & F^T = I \\ F^T & F = I \end{matrix}$$
- Technical:
  - Use orthonormal transformation matrix  $F \in \mathbb{R}^{d \times k}$
  - $(x')^T = x^T \cdot F$  is the transformation of (column) vector  $x$  into the new coordinate system defined by  $F$
  - $X' = X \cdot F$  is the matrix containing all the transformed points  $x'_i$

$$f [ ] d$$

A diagram illustrating matrix multiplication. On the left is a vertical rectangle labeled  $N \times d$ . To its right is a small square labeled  $d \times k$ . To the right of that is another vertical rectangle labeled  $N \times k$ . Below the first rectangle is the letter  $X$ , below the second is  $F$ , and below the third is  $X'$ . Between the first and second rectangles is a multiplication sign. Between the second and third rectangles is an equals sign.



# Discussion: Linear Transformations

- Feature selection is a linear transformation
  - What is matrix  $F$ ?  
$$F = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}$$

2

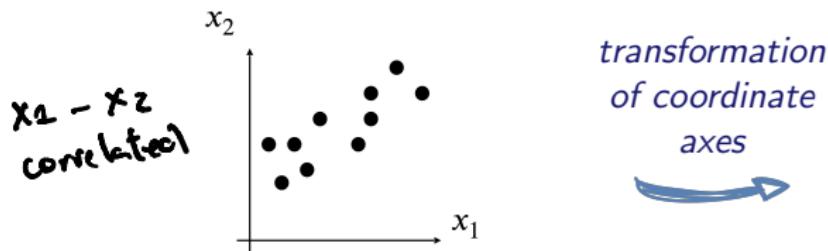
*discards 2nd dimension*
  - Let  $\bar{x}$  be the mean vector (here: row vector) in the original data space, the mean vector in the transformed space is given by  
$$\bar{x}' = \bar{x} \cdot F$$
    - Proof?
  - Let  $\Sigma_X$  be the covariance matrix in the original data space, the covariance matrix in the transformed space is then  
$$\Sigma_{X'} = F^T \cdot \Sigma_X \cdot F$$
    - Proof?

## Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. **Principal Component Analysis (PCA)**
3. Singular Value Decomposition (SVD)
4. Matrix Factorization
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

# Principal Component Analysis: Motivation

- Question: Which transformation matrix  $F$  to use?
  - Is there an **optimal orthogonal transformation** (depending on the data)?
  - Optimality: Approximate the data with few coefficients as well as possible



- Approach: Principal Component Analysis (PCA)
  - Find a coordinate system in which the (possibly originally correlated) points are **linearly uncorrelated**
  - The dimensions with no or low variance can then be ignored

# Determine the Principal Components

$$\text{Cov} : ; = 0$$

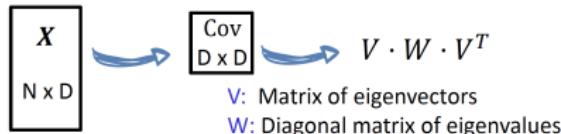
- Goal:
  - Transform the data, such that the **covariance between the new dimensions is 0**
  - The transformed data points are not linearly correlated any more
- Given:  $N$   $d$ -dimensional data points:  
 $\{\mathbf{x}_i\}_{i=1}^N, \mathbf{x}_i \in \mathbb{R}^d \quad \forall i \in \{1, \dots, N\}$
- We represent this set of points by a matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Nd} \end{bmatrix} \xrightarrow{\text{row} \rightarrow \text{instance}}$$

- The row  $\mathbf{x}_i = \{x_{i1}, \dots, x_{id}\} \in \mathbb{R}^d$  denotes the  $i$ -th point and the column  $\mathbf{X}_{:,j}$  denotes the vector containing all values from the  $j$ -th dimension

# Determine the Principal Components

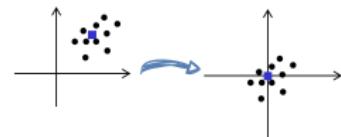
- Goal:
  - Transform the data, such that the **covariance between the new dimensions is 0**
  - The transformed data points are not linearly correlated any more
- General approach
  1. Center the data *→ zero mean*
  2. Compute the covariance matrix
  3. Use the Eigenvector Decomposition to transform the coordinate system



# Determine the Principal Components

## 1. Centering the data

- Given:  $X \in \mathbb{R}^{N \times d}$ :  $X = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Nd} \end{bmatrix}$
- Shift the points by their mean  $\bar{x} \in \mathbb{R}^d$   
(centralized data):  $\tilde{x}_i = x_i - \bar{x}$



Statistics:

$$\bar{x}_i = \frac{1}{N} \sum x_i$$

**Zero order statistic:** number of points  $N$

**First order statistic:** mean of the  $N$  points, the vector  $\bar{x} \in \mathbb{R}^d$ :

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_d \end{bmatrix} = \frac{1}{N} \cdot X^T \cdot \mathbf{1}_N \quad (1)$$

*Zentral*

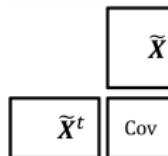
where  $\mathbf{1}_N$  is an  $N$ -dimensional vector of ones

## Determine the Principal Components

### 2. Compute the covariance Matrix

$$\Sigma_{\tilde{\mathbf{X}}} = \begin{bmatrix} \text{Var}(\tilde{\mathbf{X}}_1) & \dots \\ \dots & \text{Cov}(\tilde{\mathbf{X}}_{i,i}) \end{bmatrix}$$

- Determine the variances  $\text{Var}(\tilde{\mathbf{X}}_j)$  for each dimension  $j \in \{1, \dots, d\}$
- Determine the covariance  $\text{Cov}(\tilde{\mathbf{X}}_{j_1}, \tilde{\mathbf{X}}_{j_2})$  between dimensions  $j_1$  and  $j_2$ ,  $\forall j_1 \neq j_2 \in \{1, \dots, d\}$
- Leads to covariance matrix  $\Sigma_{\tilde{\mathbf{X}}} \in \mathbb{R}^{d \times d}$



**Second order statistic:** variance and covariance

The variance within the  $j$ -th dimension in  $\mathbf{X}$  is:

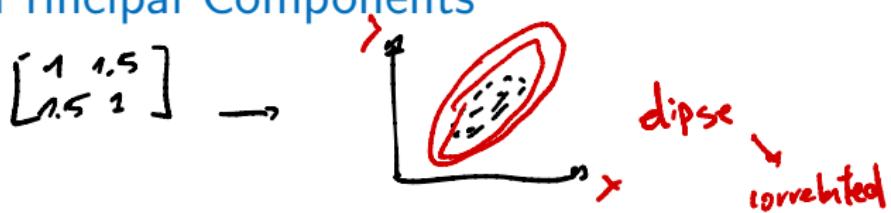
$$\text{Var}(\mathbf{X}_j) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_{ij} - \bar{\mathbf{x}}_j)^2 = \frac{1}{N} \mathbf{X}_j^T \mathbf{X}_j - \bar{\mathbf{x}}_j^2 \quad (2)$$

The covariance between dimensions  $j_1$  and  $j_2$  is:

$$\mathbf{x} = \begin{bmatrix} | & | \\ \mathbf{j}_1 & \mathbf{j}_2 \end{bmatrix}$$

$$\text{Cov}(\mathbf{X}_{j_1}, \mathbf{X}_{j_2}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_{ij_1} - \bar{\mathbf{x}}_{j_1})(\mathbf{x}_{ij_2} - \bar{\mathbf{x}}_{j_2}) = \frac{1}{N} \mathbf{X}_{j_1}^T \mathbf{X}_{j_2} - \bar{\mathbf{x}}_{j_1} \bar{\mathbf{x}}_{j_2} \quad (3)$$

# Determine the Principal Components



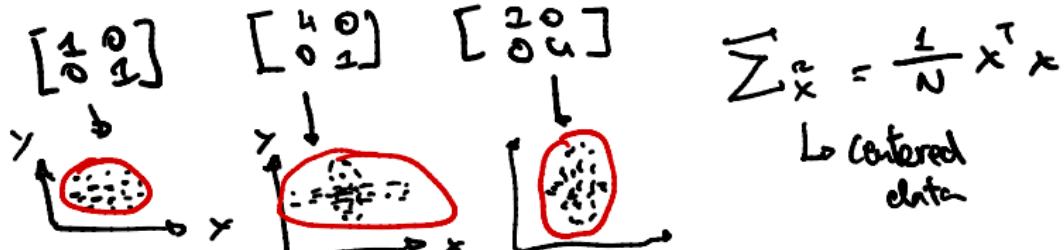
## Statistics:

For the set of points contained in  $\mathbf{X}$  the corresponding covariance matrix is defined as:

$$\Sigma_{\mathbf{X}} = \begin{bmatrix} \text{Var}(\mathbf{x}_1) & \text{Cov}(\mathbf{x}_1, \mathbf{x}_2) & \dots & \text{Cov}(\mathbf{x}_1, \mathbf{x}_d) \\ \text{Cov}(\mathbf{x}_2, \mathbf{x}_1) & \text{Var}(\mathbf{x}_2) & & \\ \vdots & & \ddots & \\ \text{Cov}(\mathbf{x}_d, \mathbf{x}_1) & \dots & & \text{Var}(\mathbf{x}_d) \end{bmatrix} = \frac{1}{N} \mathbf{X}^T \mathbf{X} - \bar{\mathbf{x}} \bar{\mathbf{x}}^T \quad (4)$$

~ zero

Remark: Covariance matrices are symmetric



# Determine the Principal Components

**Goal of PCA:** Transformation of the coordinate system such that the covariances between the new axes are 0

- Approach
  - Diagonalization by changing the basis (= adapt the coordinate system)
  - According to the spectral theorem, the eigenvectors of a symmetric matrix form an orthogonal basis
- Eigendecomposition of the covariance matrix:  $\Sigma_{\tilde{X}} = \Gamma \cdot \Lambda \cdot \Gamma^T$

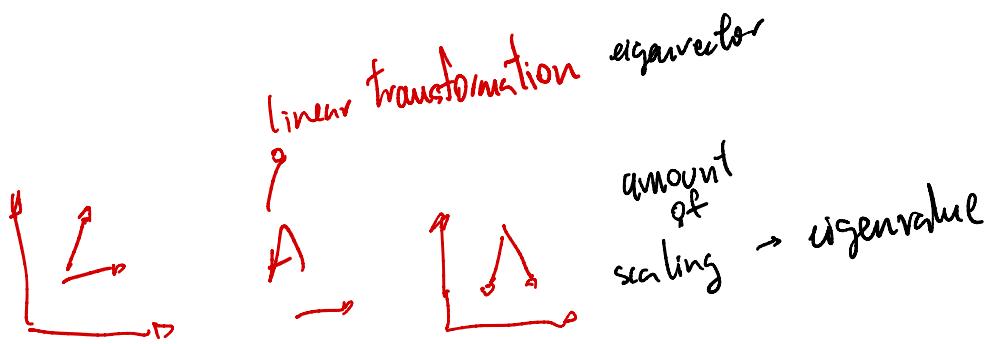
$$\Lambda = \begin{bmatrix} Var(1)' & 0 & \dots & 0 \\ 0 & Var(2)' & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & Var(D)' \end{bmatrix}$$

diagonal matrix  $\hookrightarrow \Lambda = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \rightarrow$  eigenvalues of the covariance matrix  
gamma under  $\Gamma$  gamma transpose

Eigendecomposition (spectral decomposition) is the factorization of  $X \in \mathbb{R}^{d \times d}$ :

$$X = \Gamma \cdot \Lambda \cdot \Gamma^T \quad (5)$$

- Matrices  $\Gamma, \Lambda \in \mathbb{R}^{d \times d}$  with columns of  $\Gamma$  being the normalized eigenvectors  $\gamma_i$
- $\Gamma$  is an orthonormal matrix:  $\Gamma \cdot \Gamma^T = \Gamma^T \cdot \Gamma = Id (\Gamma^T = \Gamma^{-1})$
- $\Lambda$  is a diagonal matrix with eigenvalues  $\lambda_i$  as the diagonal elements



$$Ar = \lambda r$$

# Determine the Principal Components

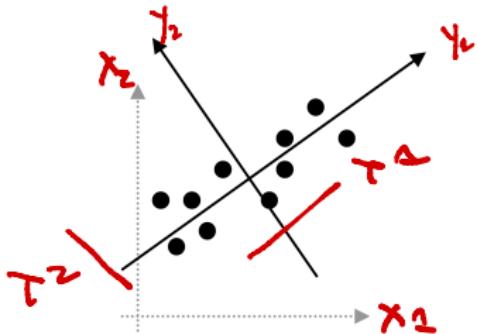
- Eigendecomposition of the covariance matrix:  $\Sigma_{\tilde{X}} = \Gamma \cdot \Lambda \cdot \Gamma^T$

$$\Sigma_{\tilde{X}} = \Gamma \sum_{\tilde{X}} \Gamma^T = \underbrace{\Gamma^T \Gamma}_{I} \Lambda \underbrace{\Gamma^T \Gamma}_{I} = \Lambda = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \rightarrow \text{eigen values}$$

(slide 9)

- The **new coordinate system** is defined by the eigenvectors  $\gamma_i$ :

- Transformed data:  $\mathbf{Y} = \tilde{\mathbf{X}} \cdot \Gamma$
- $\Lambda$  is the covariance matrix in this new coordinate system
  - New system has variance  $\lambda_i$  in dimension  $i$
  - $\forall i_1 \neq i_2 : \text{Cov}(\mathbf{Y}_{i_1}, \mathbf{Y}_{i_2}) = 0$



# Dimensionality Reduction with PCA

$$T = \begin{bmatrix} | & | \end{bmatrix} \quad \Lambda = \begin{bmatrix} 10 & & \\ & 2 & \\ & & 1 & \dots \end{bmatrix}$$

- Approach
  - The coordinates with low variance (hence low  $\lambda_i$ ) can be ignored
  - W.l.o.g. let us assume  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

- Truncation of  $\Gamma$ 
  - Keep only columns (i.e. eigenvectors) of  $\Gamma$  corresponding to the largest  $k$  eigenvalues  $\lambda_1, \dots, \lambda_k$  and ignore the rest
  - $Y_{\text{reduced}} = \tilde{X} \cdot \Gamma_{\text{truncated}}$

- How to pick  $k$ ?
  - Frequently used: 90% rule; the  $k$  variances should explain 90% of the energy
  - $k = \text{smallest value ensuring } \sum_{i=1}^k \lambda_i \geq 0.9 \cdot \sum_{i=1}^d \lambda_i$

- The modified points (transformed and truncated) contain most of the information of the original points and are low dimensional

$$\tilde{x} = x - \bar{x}$$
$$\bar{x} \Gamma = \tilde{y}$$

# Complexity

- Complexity of PCA:

$$O(N \cdot d^2) + (d^3) + O(N \cdot d \cdot k) = O(N \cdot d^2 + d^3)$$

Compute              Eigenvalue              Project data onto  
covariance matrix    decomposition    the k-dimensional space

- Remarks on eigenvalue decomposition:

- Usually we are interested in the reduced data only
- **Only the k largest eigenvectors required** (i.e. not all of them)
- Use iterative approaches (next slide) for finding eigenvectors
  - Complexity:  $O(\#it \cdot d^2)$  // #it= number of iterations
  - For sparse data even faster:  $O(\#it \cdot \#nz)$  // #nz= number of nonzero elements in the matrix

# How to Compute Eigenvectors?

- Eigenvalues are important for many machine learning/data mining tasks
  - PCA, Ranking of Websites, Community Detection, . . . //see our other lecture!
  - How to compute them efficiently?
- Power iteration (a.k.a. Von Misesiteration)
  - Iterative approach to compute **a single** eigenvector
- Let  $A$  be a matrix and  $v$  be an arbitrary (normalized) vector
  - Iteratively compute  $v \leftarrow \frac{A \cdot v}{\|A \cdot v\|}$ 
    - In each step,  $v$  is simply multiplied with  $A$  and normalized
    - **$v$  converges to the eigenvector of  $A$  with largest absolute value**
    - Highly efficient for sparse data

$$A = \sum \tilde{x}$$

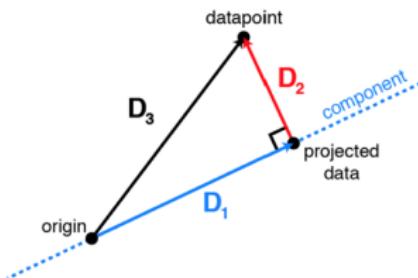
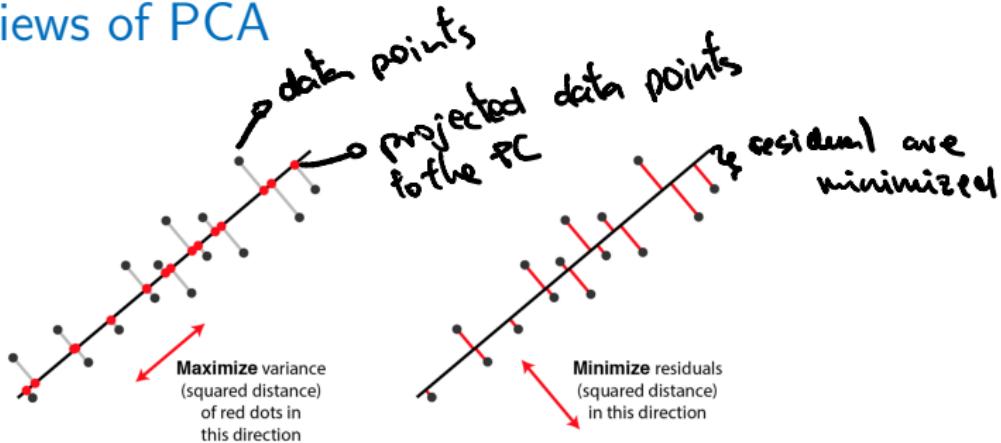
# How to Compute Eigenvectors?

depends on how quick  
it takes to converge

- Converges
  - Linear convergence with rate  $|\delta_2/\delta_1|$
  - Fast convergence if first and second eigenvalue are dissimilar
- How to find **multiple (the k largest) eigenvectors?**
  - Let us focus on symmetric matrices  $A$
  - Eigenvalue decomposition leads to:  $A = \Gamma \cdot \Lambda \cdot \Gamma^T = \sum_{i=1}^d \lambda_i \cdot \gamma_i \cdot \gamma_i^T$
  - Define deflated matrix:  $\tilde{A} = A - \lambda_1 \cdot \gamma_1 \cdot \gamma_1^T$ 
    - $\tilde{A}$  has the same eigenvectors as  $A$  except the first one has become zero
  - Apply power iteration on  $\tilde{A}$  to find the second largest eigenvector of  $A$

$$O(d^2 \cdot \text{fit} \cdot k) \rightarrow \text{total complexity}$$

# Alternative views of PCA



$$D_3^2 = D_1^2 + D_2^2$$

$\downarrow$  min       $\downarrow$  min

Images adapted from Alex H. Williams

# PCA vs. Regression

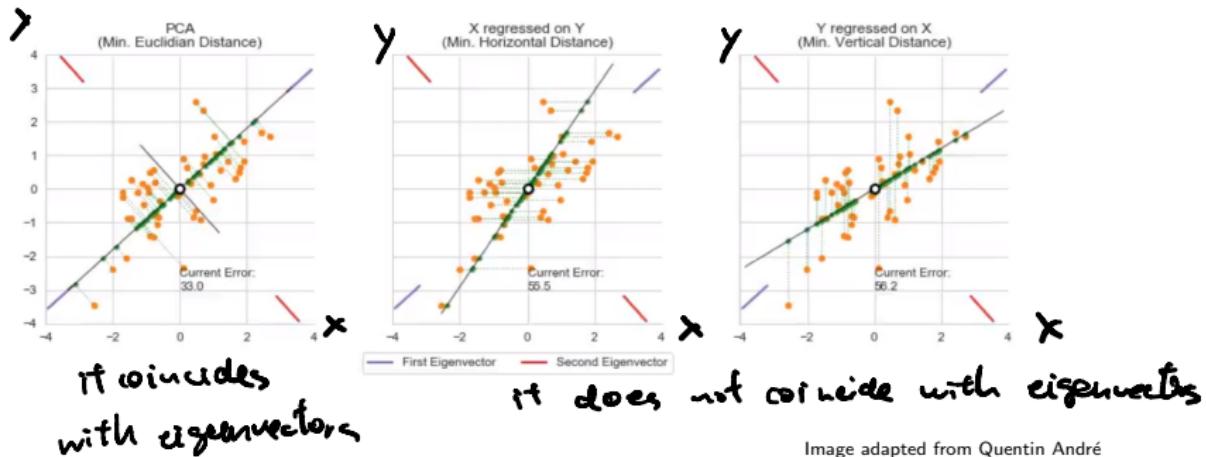


Image adapted from Quentin André

# PCA: Summary

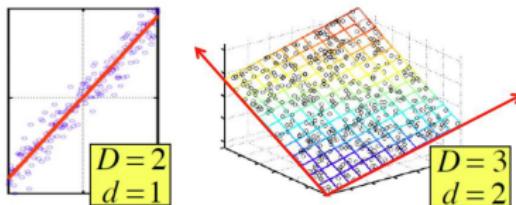
- PCA finds the optimal transformation by deriving uncorrelated dimensions
  - Exploits eigendecomposition
- Dimensionality reduction
  - After transformation simply remove dimensions with lowest variance (or use only the  $k$  largest eigenvectors for transformation)
- Limitations
  - Only captures linear relationships (one solution: Kernel PCA)

## Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. **Singular Value Decomposition (SVD)**
  - Idea: Low Rank Approximation
  - SVD & Latent Factors
  - Dimensionality Reduction
4. Matrix Factorization
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

# Low-Dimensional Manifold

- Data often lies on a low-dimensional manifold embedded in higher dimensional space



- How can we measure the dimensionality of this manifold?
  - Put differently: how to measure the intrinsic dimensionality of the data
- How can we find this manifold?

# Rank of a Matrix

- Q: What is the rank of a matrix  $A$ ?
- A: Number of linearly independent columns/rows of  $A$
- Example:

- Matrix  $A = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$  has a rank  $r = 2$

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.

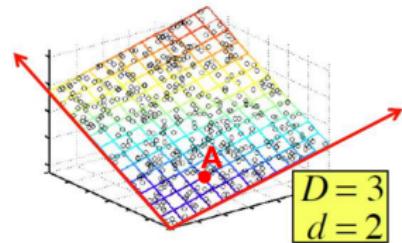
- Why do we care about low rank?
  - We can write  $A$  as two "basis" vectors:  $\begin{bmatrix} 1 & 2 & 1 \end{bmatrix}, \begin{bmatrix} -2 & -3 & 1 \end{bmatrix}$
  - And new coordinates of:  $\begin{bmatrix} 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & -1 \end{bmatrix}$

$$\begin{bmatrix} 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} -2 & -3 & 1 \end{bmatrix} \begin{bmatrix} 1+2, 2+3, 1-1 \end{bmatrix}$$

# Rank is "Dimensionality"

- Cloud of points in 3D space:
  - Think of point positions as a matrix:

1 row per point:  $\begin{pmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{pmatrix}$  A  
B  
C



- 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!

# Low Rank Approximation

- Idea: approximate original data  $\mathbf{A}$  by a low rank matrix  $\mathbf{B}$

$$\mathbf{A} = \begin{bmatrix} 1.01 & 2.05 & 0.9 \\ -2.1 & -3.05 & 1.1 \\ 2.99 & 5.01 & 0.3 \end{bmatrix} \approx \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} = \mathbf{B}$$

$$rank(\mathbf{A}) = 3$$

We need 3 coordinates  
to describe each point

$$rank(\mathbf{B}) = 2$$

We need only 2 coordinates  
per point

# Low Rank Approximation

- Idea: approximate original data  $A$  by a low rank matrix  $B$

$$A = \begin{bmatrix} 1.01 & 2.05 & 0.9 \\ -2.1 & -3.05 & 1.1 \\ 2.99 & 5.01 & 0.3 \end{bmatrix} \approx \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} = B$$

- Important: Even though both  $A$  and  $B$  are  $\in \mathbb{R}^{n \times d}$  we need only two coordinates per point to describe  $B$

- $\text{rank}(A) = 3$  vs.  $\text{rank}(B) = 2$  (3 vs. 2 coordinates per point)

- Goal: Find the best low rank approximation

- best = minimize the sum of reconstruction error
- Given matrix  $A \in \mathbb{R}^{n \times d}$ , find  $B \in \mathbb{R}^{n \times d}$  with  $\text{rank}(B) = k$  that minimizes

Find  $B$

$$\|A\|_F^2 = \sqrt{\sum_i \sum_j A_{ij}^2} \quad \min_B \quad \|A - B\|_F^2 = \sum_{i=1}^N \sum_{j=1}^D (a_{ij} - b_{ij})^2 \quad (6)$$

*rank(B) ≤ k*

## Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. **Singular Value Decomposition (SVD)**
  - Idea: Low Rank Approximation
  - **SVD & Latent Factors**
  - Dimensionality Reduction
4. Matrix Factorization
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

# Singular Value Decomposition (SVD): Definition

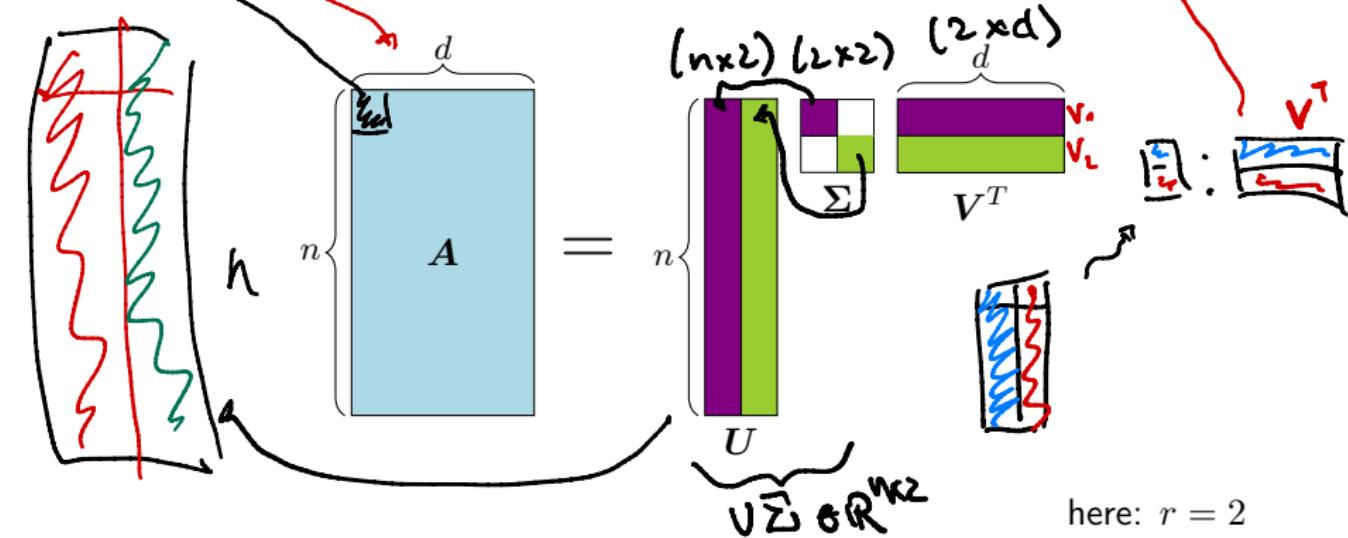


- Each real matrix  $A \in \mathbb{R}^{n \times d}$  can be decomposed into  $A = U \cdot \Sigma \cdot V^T$  (note: **exact representation**, not approximation), where
  - $U \in \mathbb{R}^{n \times r}, \Sigma \in \mathbb{R}^{r \times r}, V \in \mathbb{R}^{d \times r}$   $r \leq \min(n, d)$
  - $U, V$ : column orthonormal
    - i.e.  $UU^T = I; VV^T = I$  ( $I$  identity matrix)
    - $U$  are called the left singular vectors,  $V$  the right singular vectors
  - $\Sigma$ : diagonal
    - $r \times r$  diagonal matrix ( $r$  rank of matrix  $A$ )
    - Entries (called singular values) are positive, and sorted in decreasing order ( $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$ )
- Remark: The decomposition is (almost) unique
  - see e.g. multiplication by -1

## Singular Value Decomposition

$$h \times 2 \quad \left\{ \begin{pmatrix} 13 \\ 21 \\ n_2 \end{pmatrix} \quad \begin{pmatrix} 12 \\ 22 \end{pmatrix} \right. \\ \left. 2 \times 2 \quad \begin{pmatrix} 5 & 5 \\ 4 & 4 \\ 5 & 5 \end{pmatrix} \right.$$

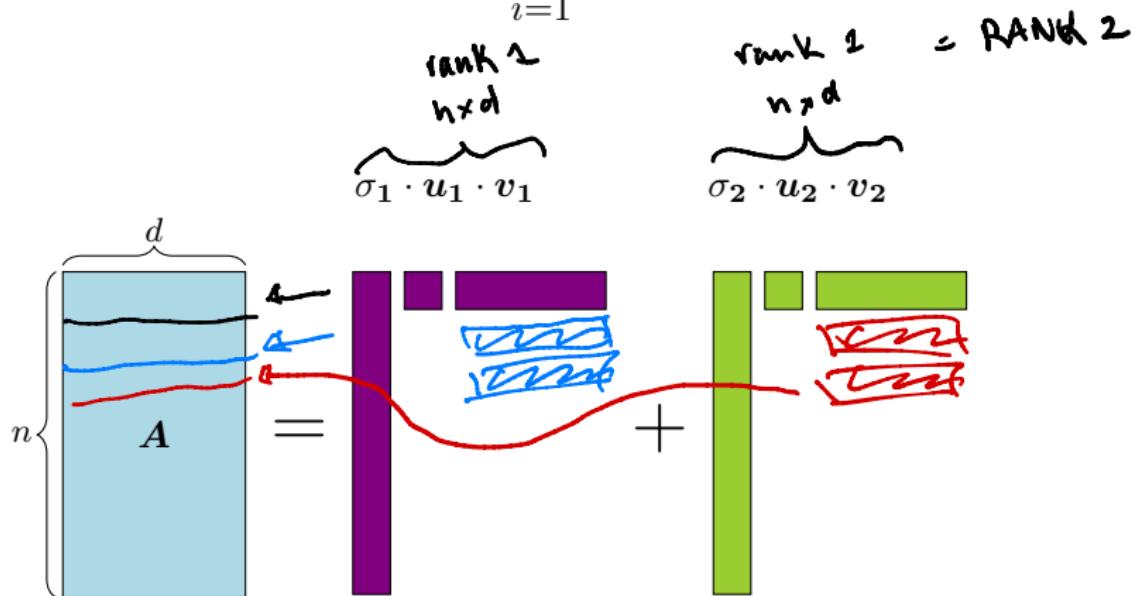
$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i \cdot u_i \cdot v_i^T \quad (7)$$



here:  $r = 2$

# Singular Value Decomposition

$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i \cdot u_i \cdot v_i^T$$



# SVD Example: Users-to-Movies

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

2

$$\text{Users Matrix} \begin{bmatrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} \\ 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix}_{(7 \times 2)} = \begin{bmatrix} 0.14 & 0 \\ 0.42 & 0 \\ 0.56 & 0 \\ 0.70 & 0 \\ 0 & 0.59 \\ 0 & 0.74 \\ 0 & 0.29 \end{bmatrix}_{(7 \times 2)} \times \Sigma \begin{bmatrix} 12.36 & 0 \\ 0 & 9.48 \end{bmatrix}_{(2 \times 2)} \times \begin{bmatrix} 0.57 & 0.57 & 0.57 & 0 & 0 \\ 0 & 0 & 0 & 0.70 & 0.70 \end{bmatrix}_{(2 \times 5)}$$

SciFi-concept      Romance-concept

Σ

$\Sigma$

$V^T$

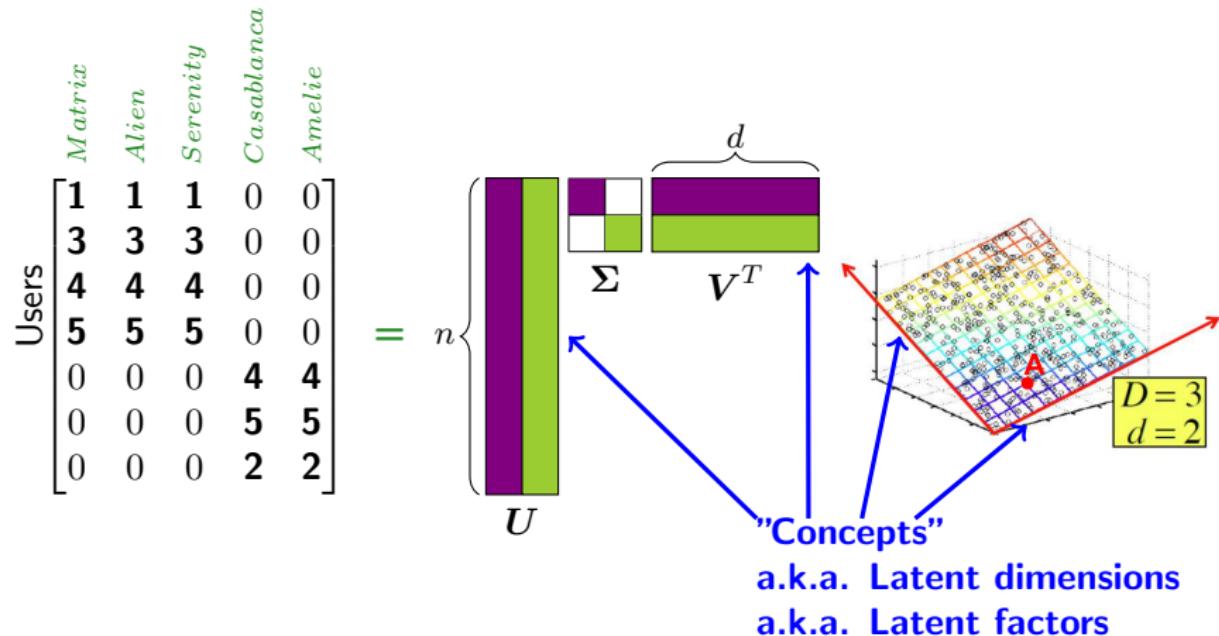
# SVD Example: Users-to-Movies

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

$$\text{Users} \begin{bmatrix} \text{Matrix} \\ 1 & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.14 & 0 \\ 0.42 & 0 \\ 0.56 & 0 \\ 0.70 & 0 \\ 0 & -0.59 \\ 0 & -0.74 \\ 0 & -0.29 \end{bmatrix} \times \begin{bmatrix} \text{SciFi-concept} & \text{Romance-concept} \\ (\text{note the multiplication by } -1) \end{bmatrix} \times \begin{bmatrix} 12.36 & 0 \\ 0 & 9.48 \\ 0.57 & 0.57 & 0.57 & 0 & -0.70 & -0.70 \end{bmatrix}$$

# SVD Example: Latent Factors

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



## SVD Example: Beyond Blocks

RANK 3

Matrix

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times$$
$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD Example: Beyond Blocks

Matrix

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

SciFi-concept      Romance-concept

$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD Example: Beyond Blocks

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

$$\begin{array}{l}
 \text{Matrix} \\
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = U \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \Sigma \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} V^T
 \end{array}$$

SciFi-concept       $U$       Romance-concept

↴      ↪      ↪

↴      ↪      ↪

↴      ↪      ↪

?

?

?

?

?

?

# SVD Example: Beyond Blocks

Matrix

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{matrix} \text{SciFi-concept} \\ \downarrow \\ \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \end{matrix} \quad \begin{matrix} \text{"strength" of the SciFi-concept} \\ \downarrow \\ \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \end{matrix}$$

**0.13**    **0.41**    **0.55**    **0.68**    **0.15**    **0.07**    **0.07**

**0.59**    **0.65**    **0.73**    **-0.67**    **0.29**    **0.32**

**12.4**    **0**    **0**

**9.5**    **0**    **0**

**1.3**    **0**    **0**

$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD Example: Beyond Blocks

$V$  is "movie-to-concept" similarity matrix

Matrix

			<i>Alien</i>	<i>Serenity</i>	<i>Casablanca</i>	<i>Amelie</i>
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		

$\downarrow$  SciFi-concept

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

SciFi-concept

$$\begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD: Interpretation

- $A = U\Sigma V^T$
- 'movies', 'users' and 'concepts':
  - $A$ : original data: movies-to-users
  - $U$ : user-to-concept similarity matrix
  - $V$ : movie-to-concept similarity matrix
  - $\Sigma$ : its diagonal elements: 'strength' of each concept
- **Benefits of SVD (or in general matrix decomposition):**
  - Discover hidden correlations/topics
    - Words that occur commonly together; movies of the same genre; ...
  - Interpretation and visualization

## Chapter: Dimensionality Reduction & Matrix Factorization

1. Introduction
2. Principal Component Analysis (PCA)
3. **Singular Value Decomposition (SVD)**
  - Idea: Low Rank Approximation
  - SVD & Latent Factors
  - **Dimensionality Reduction**
4. Matrix Factorization
5. Neighbor Graph Methods
6. Autoencoders (Nonlinear Dimensionality Reduction)

## Recap: Dim. Reduction by Low Rank Approx.

- Idea: approximate original data  $\mathbf{A}$  by a low rank matrix  $\mathbf{B}$

$$\mathbf{A} = \begin{bmatrix} 1.01 & 2.05 & 0.9 \\ -2.1 & -3.05 & 1.1 \\ 2.99 & 5.031 & 0.3 \end{bmatrix} \approx \mathbf{B} = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$$

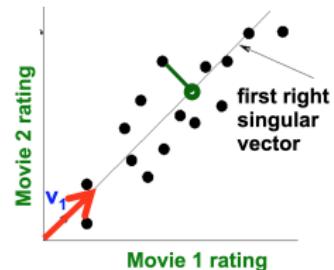
- Important: Even though both  $\mathbf{A}$  and  $\mathbf{B}$  are  $\in \mathbb{R}^{n \times d}$  we need only two coordinates per point to describe  $\mathbf{B}$ 
  - $rank(\mathbf{A}) = 3$  vs.  $rank(\mathbf{B}) = 2$  (3 vs. 2 coordinates per point)
- Goal: Find the best low rank approximation**
  - best = minimize the sum of reconstruction error
  - Given matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , find  $\mathbf{B} \in \mathbb{R}^{n \times d}$  with  $rank(\mathbf{B}) = k$  that minimizes

$$\|\mathbf{A} - \mathbf{B}\|_F^2 = \sum_{i=1}^N \sum_{j=1}^D (a_{ij} - b_{ij})^2 \quad (8)$$

# SVD: Alternative Interpretation

*Singular values*

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

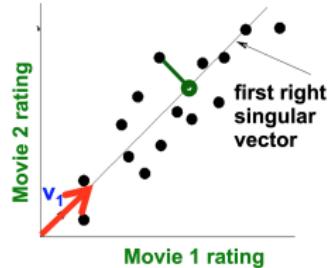


Variance ("spread")  
on the  $v_1$  axis

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD: Alternative Interpretation

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



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

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: Best Approximation

- How to find the best approximation?

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} \mathbf{0.13} & -0.02 & -0.01 \\ \mathbf{0.41} & -0.07 & -0.03 \\ \mathbf{0.55} & -0.09 & -0.04 \\ \mathbf{0.68} & -0.11 & -0.05 \\ 0.15 & \mathbf{0.59} & \mathbf{0.65} \\ 0.07 & \mathbf{0.73} & \mathbf{-0.67} \\ 0.07 & \mathbf{0.29} & \mathbf{0.32} \end{bmatrix} \times \begin{bmatrix} \mathbf{12.4} & 0 & 0 \\ 0 & \mathbf{9.5} & 0 \\ 0 & 0 & \mathbf{1.3} \end{bmatrix} \times \begin{bmatrix} \mathbf{0.56} & \mathbf{0.59} & \mathbf{0.56} & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & \mathbf{0.69} & \mathbf{0.69} \\ 0.42 & \mathbf{-0.80} & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD: Best Approximation

- How to find the best approximation?
- Set smallest singular values to zero!

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD: Best Approximation

- How to find the best approximation?
- Set smallest singular values to zero!

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 13 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.42 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

# SVD: Best Approximation

- How to find the best approximation?
- Set smallest singular values to zero!

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

# SVD: Best Approximation

- How to find the best approximation?
- Set smallest singular values to zero!

By construction,  
the rank of the  
new matrix is 2

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

# SVD: Best Low Rank Approximation

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

**$B$  is best approximation of  $A$**

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

## SVD: Projection

- Note: The actual projected/reduced data can be obtained by computing

$$P = U \Sigma$$

- Or equivalently:  $P = A \cdot V$  (since  $V$  is orthonormal)

# Best Approximation – Intuitive Explanation

- Recap: Vectors  $u_i$  and  $v_i$  are of unit length
- W.l.o.g.:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$

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

# Best Approximation – Intuitive Explanation

- Recap: Vectors  $u_i$  and  $v_i$  are of unit length
- W.l.o.g.:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \underbrace{\sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots}_{r \text{ terms}}$$

**Q: How many  $\sigma_i$  to pick?**

**A:** Rule of thumb:

keep 90% of 'energy'

$$\sum_{i=1}^k \sigma_i^2 = 0.9 \sum_{i=1}^r \sigma_i^2$$

- $\sigma_i$  scales the terms  $u_i \cdot v_i^T$
- Zeroing small  $\sigma_i$  introduces less error

# SVD: Best Low Rank Approximation – Proof

- Theorem: Let  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$  ( $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$ ,  $\text{rank}(\mathbf{A}) = r$ ) and  $\mathbf{B} = \mathbf{U}\mathbf{S}\mathbf{V}^T$  with  $\mathbf{S}$  being a diagonal  $r \times r$  matrix where

- $s_i = \sigma_i$  for  $i = 1 \dots k$  and  $s_i = 0$

Then  $\mathbf{B}$  is a best rank=k approximation to  $\mathbf{A}$  regarding Frobenius norm, i.e.  $\mathbf{B}$  is a solution to  $\min_{\mathbf{B}} \|\mathbf{A} - \mathbf{B}\|_F$  where  $\text{rank}(\mathbf{B}) = k$

- We have uploaded a detailed proof to the web
  - Note: Many proofs on the web and on other lecture slides are incorrect!
- Remark:  $\mathbf{B}$  is also an optimal low-rank approximation regarding the spectral norm (operator 2-norm):  $\min_{\mathbf{B}} \|\mathbf{A} - \mathbf{B}\|_2$ 
  - $\|\mathbf{X}\|_2 = \text{largest singular value of } \mathbf{X}$

# SVD: Best Low Rank Approximation – Proof

Some facts:

- $\|\mathbf{X}\|_F = \|\mathbf{X}^T\|_F$ 
  - obvious from the definition
- $\|\mathbf{X}\|_F^2 = \text{trace}(\mathbf{X}^T \mathbf{X})$
- Frobenius norm is invariant to orthonormal transformations  $\mathbf{U}$ 
  - Note: if  $\mathbf{U}^T \mathbf{U} = I$  then also  $\mathbf{U} \mathbf{U}^T = I$
  - $\|\mathbf{U} \mathbf{X}\|_F^2 = \text{trace}((\mathbf{U} \mathbf{X})^T (\mathbf{U} \mathbf{X})) = \text{trace}(\mathbf{X}^T \mathbf{U}^T \mathbf{U} \mathbf{X}) = \text{trace}(\mathbf{X}^T \mathbf{X}) = \|\mathbf{X}\|_F^2$
  - $\|\mathbf{X} \mathbf{U}\|_F^2 = \|(\mathbf{X} \mathbf{U})^T\|_F^2 = \text{trace}((\mathbf{X} \mathbf{U})(\mathbf{X} \mathbf{U})^T) = \text{trace}(\mathbf{X} \mathbf{U} \mathbf{U}^T \mathbf{X}^T) = \text{trace}(\mathbf{X} \mathbf{X}^T) = \|\mathbf{X}\|_F^2$
- Let  $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$  then  $\|\mathbf{A}\|_F^2 = \|\Sigma\|_F^2 = \sum_i^r \sigma_i^2$ 
  - follows from above results

# SVD: Complexity

- To compute SVD:
  - $O(n \cdot d^2)$  or  $O(n^2 \cdot d)$  (whichever is smaller)
- 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 & PCA: Comparison

- Given data  $\mathbf{X}$  (let's assume it is already centered)
- SVD gives us
  - $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T$
  - Projected data obtained by  $\mathbf{X}\mathbf{V} = \mathbf{U}\Sigma$  (or truncated  $\mathbf{V}$ )
- PCA computes the eigendecomposition of the covariance matrix
  - Covariance matrix:  $\mathbf{X}^T\mathbf{X}$
  - Eigendecomposition leads to  $\mathbf{X}^T\mathbf{X} = \mathbf{\Gamma} \cdot \mathbf{\Lambda} \cdot \mathbf{\Gamma}^T$
  - Projected data obtained by  $\mathbf{X} \cdot \mathbf{\Gamma}$  (or truncated  $\mathbf{\Gamma}$ )
- Let us calculate:
  - $\mathbf{X}^T\mathbf{X} = (\mathbf{U}\Sigma\mathbf{V}^T)^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^T\mathbf{U}^T(\mathbf{U}\Sigma\mathbf{V}^T) = \mathbf{V}\Sigma\Sigma^T\mathbf{V}^T = \mathbf{V}\Sigma^2\mathbf{V}^T$
  - $\mathbf{\Gamma} = \mathbf{V}$  **PCA and SVD are equivalent!**
  - $\Sigma^2 = \mathbf{\Lambda}$  **Squared singular values are variances in new space!**

# SVD & PCA: Comparison

Transform the data such that dimensions of new space are uncorrelated  
+ discard (new) dimensions with smallest variance

=

find optimal low-rank approximation  
(regarding Frobeniusnorm)

## Remark: Computation of SVD

- We can use the eigendecomposition to calculate the singular value decomposition
- $\mathbf{X}^T \mathbf{X} = (\mathbf{U} \Sigma \mathbf{V}^T)^T \mathbf{U} \Sigma \mathbf{V}^T = \mathbf{V} \Sigma^T \mathbf{U}^T (\mathbf{U} \Sigma \mathbf{V}^T) = \mathbf{V} \Sigma \Sigma^T \mathbf{V}^T = \mathbf{V} \Sigma^2 \mathbf{V}^T$ 
  - $\mathbf{V}$  = eigenvectors of  $\mathbf{X}^T \mathbf{X}$
- $\mathbf{X} \mathbf{X}^T = \mathbf{U} \Sigma \mathbf{V}^T (\mathbf{V} \Sigma^T \mathbf{U}^T) = \mathbf{U} \Sigma \Sigma^T \mathbf{U}^T$ 
  - $\mathbf{U}$  = eigenvectors of  $\mathbf{X} \mathbf{X}^T$
- Drawback: Numerically unstable
  - Better to use specialized algorithms