# Machine Learning

## Lecture 10: Dimensionality Reduction & Matrix Factorization

Prof. Dr. Stephan Günnemann
Aleksandar Bojchevski

Data Analytics and Machine Learning Group
Technical University of Munich

Winter term 2020/2021

# Roadmap

- 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 (Non-linear Dimensionality Reduction)

# Introduction: Unsupervised Learning (I)

- Supervised learning aims to map inputs to targets with $y = f(\boldsymbol{x})$, or in a probabilistic framework it models $p(y|\boldsymbol{x})$

- Unsupervised learning can be seen as modelling $p(\boldsymbol{x})$

  *features*

- We are trying to find the (hidden / latent) structure in the data

  - e.g. find a latent distribution $p(\boldsymbol{z})$ and a generative transformation $p(\boldsymbol{x}\,|\,\boldsymbol{z})$
    we can then obtain $p(\boldsymbol{x}) = \int p(\boldsymbol{x}\,|\,\boldsymbol{z})\,p(\boldsymbol{z})\,d\boldsymbol{z}$

    *for each $i$*

    $z_i \sim p(z)$

  - latent $\boldsymbol{z}$ usually unknown and has to be estimated

    $x_i \sim p(x|z_i)$

- Examples:

  *cluster indicator*

  - Clustering: the cluster label is the latent state

  - Anomaly detection: treat instances with low $p(\boldsymbol{x})$ as anomalies

# Introduction: Unsupervised Learning (II)

- ==Unsupervised learning== can be <u>viewed as compression</u>
  - compress a data point to a <u>single label</u> corresponding to its cluster
  - compress a data point <u>from a higher dim.</u> <u>to a lower dim.</u> <u>latent space</u>

- Unsupervised learning can be used …
  - … as a <u>stand-alone method</u> (e.g. to understand your data, visualization)
  - … as a <u>pre-processing step</u> (e.g. use cluster label as feature for subsequent classification task; obtain small number of relevant features)
  - … to <u>leverage</u> large amounts of unlabeled data for pretraining

- This lecture: Dimensionality Reduction & Matrix Factorization

# Dimensionality Reduction: Motivation

- Often data has <u>very many features</u>, i.e. high-dimensional data
- High-dimensional data is challenging:
  - Similarity search/<u>computation is expensive</u> because of high complexity of distance functions
  - <u>Highly correlated dimension</u> could cause trouble for some algorithms
  - <u>Curse of dimensionality</u>: we need exponential amounts of data to characterize the density as the dimensionality goes up
  - It is <u>hard to visualize</u> high-dimensional data
- <u>Often</u> the <u>data lies on a low-dim. manifold</u>, embedded in a high-dim. space

- Goal: Try to <u>reduce the dimen</u>sionality <u>while avoiding information loss</u>
- 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
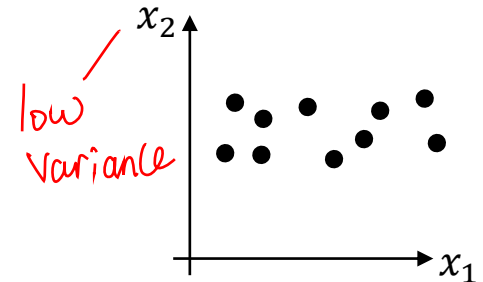    - Example: shoe size interesting for shoe purchases, not so for car purchases

- Advantages:
    - No need for an intensive preprocessing or training phase to determine relevant dimensions
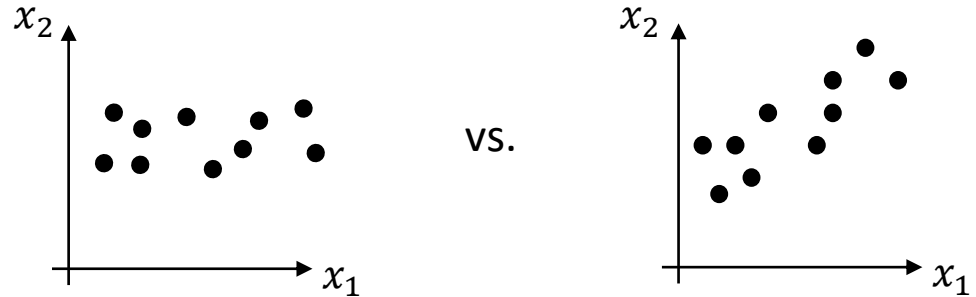
- Disadvantages:
    - Expert knowledge required; misjudgment possible
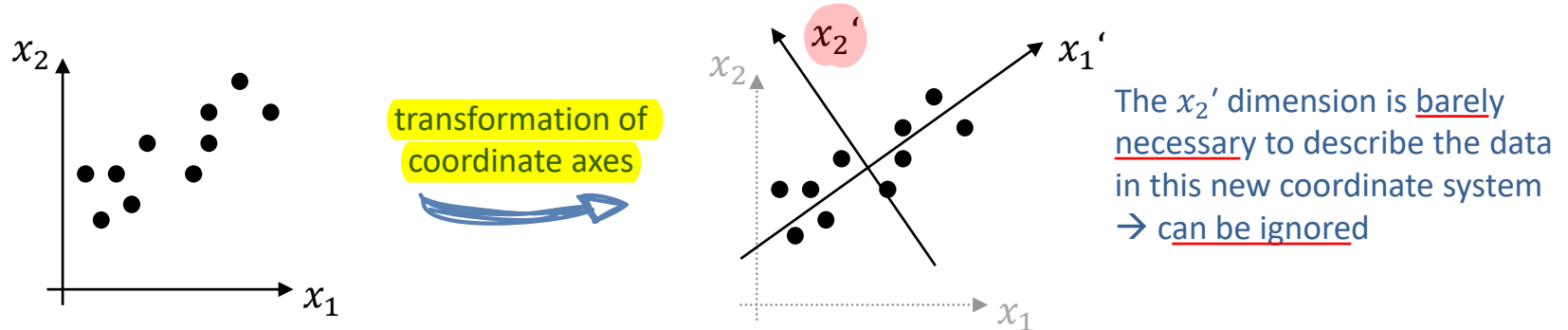    - Univariate feature selection ignores correlations

# Beyond Feature (Sub-)Selection

- Can we do
  - better?
  - automatic?



vs.

- Obviously: Simply <u>discarding whole features not a good idea</u>
  - <u>Features are often correlated</u>



transformation of coordinate axes

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

✳ Ideally, we seek to capture data independent of coord sys

# Dim. Reduction via Linear Transformations

*preserving dot product* — *length, angle*

✗ • Represent data in a different coordinate system via linear transformations

- **change of basis** (orthogonal basis transformations)

  *If such low Var dim is Known in advance, just remove them earlier*

- **+ potentially discarding dimensions** *of low variance*

- Technical:

  - use orthonormal transformation matrix $F \in R^{d \times k}$    *K < d*

  - $(x')^T = x^T \cdot F$ 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 $x'_i$

$N \times d$    * $d \times k$  = $N \times k$

$F$

$X$        $X'$

$FF^T = F^TF$

$= I$

(first center to mean, then) transformation with F

8

# Discussion: Linear Transformations

discard 2nd dim,
keep the rest

$$\text{let } F = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}_{3 \times 2}$$

- ~~Feature selection is a linear transformation~~
  - What is the matrix $F$?

- Let $\bar{x}$ be the <mark>mean vector</mark> (here: row vector) <u>in the original data space</u>, the mean vector in the transformed space is given by $\bar{x'} = \bar{x} \cdot F$

- Let $\Sigma_X$ be the <mark>covariance matrix</mark> in the original data space, the covariance matrix in the transformed space is then $\Sigma_{X'} = F^T \cdot \Sigma_X \cdot F$

# Roadmap

- 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 (Non-linear Dimensionality Reduction)

Data Analytics
and Machine Learning

# 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



Correlated

transformation of
coordinate axes

$X_2'$   $X_1'$

uncorrelated

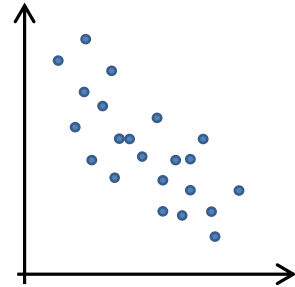- Approach: Principal Component Analysis (PCA)

  - 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

**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: $\{x_i\}_{i=1}^N$, $x_i \in \mathbb{R}^d$ $\forall i \in \{1, \ldots, N\}$

- We represent this set of points by a matrix $X \in \mathbb{R}^{N \times d}$:

*Feature matrix*

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{Nd} \end{bmatrix}$$
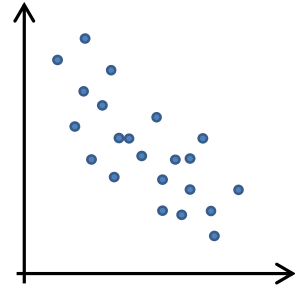
*instance*

- The row $x_i = \{x_{i1}, \ldots, x_{id}\} \in \mathbb{R}^d$ denotes the $i$-th point and the column $X_{:,j}$ denotes the vector containing all values from the $j$-th dimension
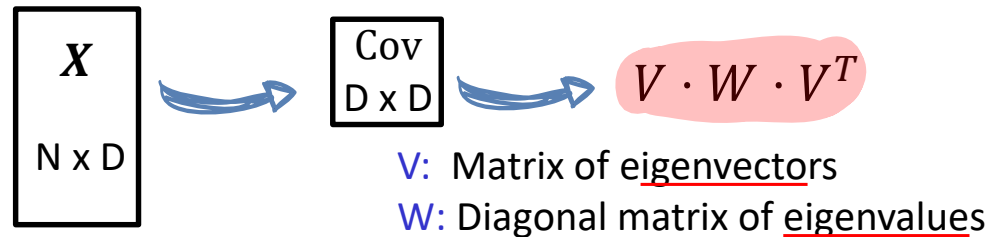
*feature*

# Determine the Principal Components

- Goal:

  - Transform the data, such that the **covariance between the new dimensions is 0**

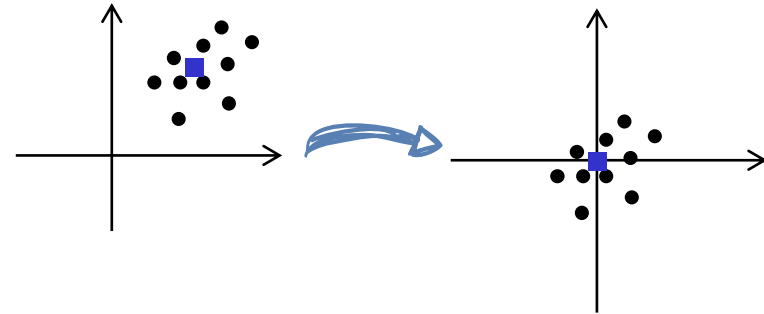  - The transformed data points are not linearly correlated any more

- General approach

  1. Center the data   *i.e Zero mean*

  2. Compute the covariance matrix

  3. Use the Eigenvector decomposition to transform the coordinate system

$$X \quad \rightarrow \quad \underset{\text{D x D}}{\text{Cov}} \quad \rightarrow \quad V \cdot W \cdot V^T$$

N x D

V: Matrix of eigenvectors
W: Diagonal matrix of eigenvalues

# Determine the Principal Components

- <u>Given</u>: $\boldsymbol{X} \in \mathbb{R}^{N \times d}$: $\boldsymbol{X} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{Nd} \end{bmatrix}$

- <u>Shift the points by their mean</u> $\overline{\boldsymbol{x}} \in \mathbb{R}^d$ (centralized data): $\widetilde{x}_i = x_i - \overline{x}$

---

<u>Statistics</u>:

Zero order statistic : number of points $N$

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

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

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

# Determine the Principal Components

- Determine the variances $\text{Var}(\widetilde{X}_j)$ for each dimension $j \in \{1, \dots d\}$

- Determine the covariance $\text{Cov}(\widetilde{X}_{j_1}, \widetilde{X}_{j_2})$ between dimensions $j_1$ and $j_2$, $\forall\, j_1 \neq j_2 \in \{1, \dots d\}$

➢ Leads to the covariance matrix $\Sigma_{\widetilde{X}} \in \mathbb{R}^{d \times d}$

$\widetilde{X}$

$\widetilde{X}^t$    Cov

Statistics:

Second order statistic: variance and covariance

The variance within the j-th dimension in $X$ is:

$$\text{Var}(X_j) = \frac{1}{N} \sum_{i=1}^{N} (x_{ij} - \overline{x}_j)^2 = \frac{1}{N} \cdot X_j^T X_j - \overline{x}_j^2$$

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

$$\text{Cov}(X_{j_1}, X_{j_2}) = \frac{1}{N} \sum_{i=1}^{N} (x_{ij_1} - \overline{x}_{j_1}) \cdot (x_{ij_2} - \overline{x}_{j_2}) = \frac{1}{N} \cdot X_{j_1}^T X_{j_2} - \overline{x}_{j_1} \overline{x}_{j_2}$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Variation in X
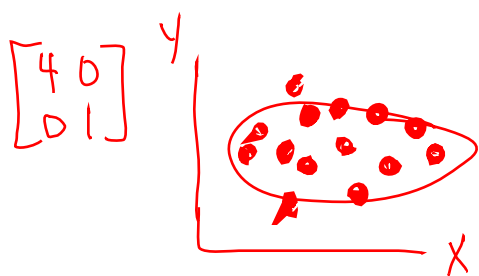equal $\Rightarrow$ Zero Correlation
Variation in Y

---

Statistics (continued):

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

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

zero mean

- Remark: Covariance matrices are symmetric

$$\begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$$

$\text{Var}_X = 4\,\text{Var}_Y$

$$\begin{bmatrix} 1 & 1.5 \\ 1.5 & 1 \end{bmatrix}$$

X & Y directly correlate
non-zero off-diag

*Our original data is usually correlated i.e non Zero values everywhere*

$$Cov' = \begin{pmatrix} Var(1)' & 0 & \cdots & 0 \\ 0 & Var(2)' & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & Var(D)' \end{pmatrix}$$

**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: $\Sigma_{\tilde{X}} = \Gamma \cdot \Lambda \cdot \Gamma^T$

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

$$A = \Gamma \cdot \Lambda \cdot \Gamma^T$$

*cols are γ$_i$ of covariance matrix*

→ matrices $\Gamma, \Lambda \in \mathbb{R}^{d \times d}$ with columns of $\Gamma$ being the normalized eigenvectors $\gamma_i$

→ $\Gamma$ is an orthonormal matrix: $\Gamma \cdot \Gamma^T = \Gamma^T \cdot \Gamma = \mathrm{Id}$ ($\Gamma^T = \Gamma^{-1}$)

→ $\Lambda$ is a diagonal matrix with eigenvalues $\lambda_i$ as the diagonal elements

*they don't need to have unit length*

$\gamma_i$ : the vec that stayed the same, only scaled by $\lambda_i$
when applying matrix A to some vec
placing $\gamma_i$ as cols for some matrix M $\Rightarrow$ M will be orthogonal

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

  – Transformed data: $Y = \widetilde{X} \cdot \Gamma$

  – $\Lambda$ is the covariance matrix in this new coordinate system

✱ New system has variance $\lambda_i$ in dimension $i$

➤ $\forall i_1 \neq i_2 : \mathrm{Cov}(Y_{i_1}, Y_{i_2}) = 0$

no correlation
bet dim

# Dimensionality Reduction with PCA ✗

- **Approach**

  - The coordinates with low variance (hence low $\lambda_i$ ) can be ignored

  - W.l.o.g. let us assume $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

- ➤ **Truncation of $\boldsymbol{\Gamma}$**

  - ✗ **Keep only columns (i.e. eigenvectors) of $\boldsymbol{\Gamma}$ corresponding to the largest $\mathbf{k}$ eigenvalues $\lambda_1, \dots, \lambda_k$**

  - $Y_{\text{reduced}} = \widetilde{X} \cdot \boldsymbol{\Gamma}_{\text{truncated}}$

  *applying linear transformation to get the new mean*

- **How to pick $k$?**

  - Frequently used: 90% rule; the *k* variances should explain 90% of the energy

  - k = 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

*GOAL*

*✗ So far it's zero mean data → Sometimes we reserve this;*
*seeking more info if lost due to subtracting the mean*

# Complexity

- Complexity of PCA:

$$\mathrm{O}(N \cdot d^2) \quad + \quad O(d^3) \quad + \quad O(N \cdot d \cdot k) = O(N \cdot d^2 + d^3)$$

<table>
<tr><td>Compute<br>covariance matrix</td><td>Eigenvalue<br>decomposition</td><td>Project data onto<br>the k-dimensional space</td></tr>
</table>

Remarks on eigenvalue decomposition:

*directly compute K largest, no need to compute all*

- 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: $O(\#\mathrm{it} \cdot d^2)$                     // #it = number of iterations

  - For sparse data even faster: $O(\#\mathrm{it} \cdot \#\mathrm{nz})$              // #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 Mises** iteration)
  - Iterative approach to compute **a single** eigenvector
  - Let $A$ be a <u>matrix</u> and $v$ be an arbitrary (<u>normalized</u>) vector
    - Iteratively compute $v \leftarrow \dfrac{A \cdot v}{\|A \cdot v\|}$ until convergence
      - 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

# How to Compute Eigenvectors?

- Convergence:    #iter depend on ↴

  - Linear convergence with rate $|\lambda_2/\lambda_1|$

  - Fast convergence if first and second eigenvalue are dissimilar
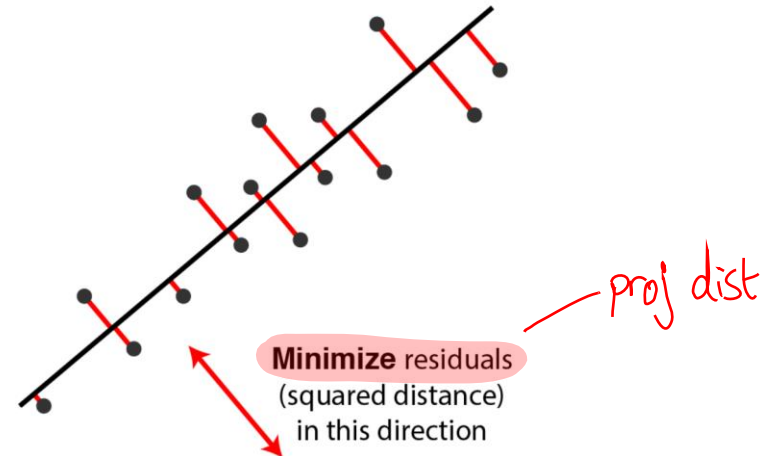
- How to find **multiple (the k largest) eigenvectors**?

  - 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: $\widehat{A} = A - \lambda_1 \cdot \gamma_1 \cdot \gamma_1^T$

    - $\widehat{A}$ has the same eigenvectors as $A$ except the first one has become zero

  - ➤ Apply power iteration on $\widehat{A}$ to find the second largest eigenvector of $A$

NB  Largest $\lambda_1$ corresponds to the dim which have the largest variance in the transformed space

# Alternative views of PCA

Eigen vector also called **Principal Component**

data
Point

proj data

**Maximize** variance
(squared distance)
of red dots in
this direction

— proj dist

**Minimize** residuals
(squared distance)
in this direction

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

— min

fixed; variance
of original space

max

datapoint

$D_3$

$D_2$

component

projected
data

origin

$D_1$

Images adapted from Alexh Williams

# PCA vs. Regression



Y given

| PCA (Min. Euclidian Distance) | X regressed on Y (Min. Horizontal Distance) | Y regressed on X (Min. Vertical Distance) |

First Eigenvector    Second Eigenvector

Current Error: 33.0

Current Error: 55.5

Current Error: 56.2

Image adapted from Quentin André

trying to min both
simultaneously

X-dir only is to
be min

$\gamma_1$ is not aligned

24

# 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)

✳ PCA is usually the first thing to be done, when analyzing data

# Roadmap

PCA is not efficient in case of non-linear data

- Chapter: Dimensionality Reduction & Matrix Factorization

    1. Introduction

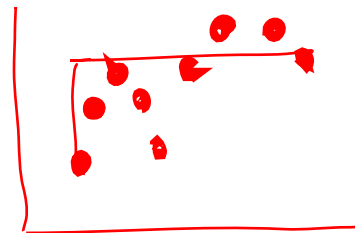    2. Principal Component Analysis (PCA) —————— equivalent

    3. **Singular Value Decomposition (SVD)**

        – **Idea: Low Rank Approximation** → aims to find the best low rank approximation for a given matrix X

        – SVD & Latent Factors

        – Dimensionality Reduction

    4. Matrix Factorization

    5. Neighbor Graph Methods

    6. Autoencoders (Non-linear Dimensionality Reduction)

# Low-Dimensional Manifold

- <u>Data</u> <u>often lies on a low</u>-<u>dimensional manifold</u> embedded in higher dimensional space



$$D = 3$$
$$d = 2$$

- How can me <u>measure the dimensionality of this manifold</u>?
  - put differently how to measure the <u>intrinsic dime</u>nsionality 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 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: [1 2 1] [-2 -3 1]
  - And new coordinates of: [1 0] [0 1] [1 -1]

# Rank is "Dimensionality"

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

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



$$D = 3$$
$$d = 2$$

**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!**

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

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

A being noisy,

its rows are no
longer dependent

all its row are
linearly independent

$\mathrm{rank}(\boldsymbol{A}) = 3$
we need 3 coordinates
to describe each point

$\mathrm{rank}(\boldsymbol{B}) = 2$
we need only 2 coordinates
per point

# Low Rank Approximation

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

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

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

  – rank(A)=3 vs. 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 $\boldsymbol{A} \in \mathbb{R}^{n \times d}$, find $\boldsymbol{B} \in \mathbb{R}^{n \times d}$ with rank($\boldsymbol{B}$) = $k$ that minimizes

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

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

We use SVD to solve it

31

# Roadmap

- 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 (Non-linear Dimensionality Reduction)

# Singular Value Decomposition (SVD): Definition

Each real matrix $A \in \mathbb{R}^{n \times d}$ can be decomposed into $A = U \cdot \Sigma \cdot V^{\mathrm{T}}$ (note: **exact representation**, no approximation), where

- $U \in \mathbb{R}^{n \times r}, \Sigma \in \mathbb{R}^{r \times r}, V \in \mathbb{R}^{d \times r}$

- $U, V$: column orthonormal
  - i.e. $U^T U = I$; $V^T V = I$ (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 $A$)     $r \leq \min(n, d)$
  - entries (called singular values) are positive, and sorted in decreasing order ($\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r > 0$)

- Remark: The decomposition is (almost) unique
  - see e.g. multiplication by -1

\* Orthonormal
- Zero dot product
- unit length

33

# Singular Value Decomposition

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



here: $r = 2$

# Singular Value Decomposition

As summation

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

Outerprod

$$\sigma_1 \cdot u_1 \circ v_1^T \qquad \sigma_2 \cdot u_2 \circ v_2^T$$

d

n $\{$ $A$ $\}$ =

rows are linearly
dependent, being
scaled by $\sigma_i$

rank = 1 by
construction

+

here: $r = 2$

35

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

**\* Assuming A is face image, U shall be called eigenfaces where rows express features**

**→ every col of A is reconstructed as linear combination of eigenfaces scaled by $\alpha_i$**

**SciFi-concept**

**Romance-concept**

$$
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
3 & 3 & 3 & 0 & 0 \\
4 & 4 & 4 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 4 & 4 \\
0 & 0 & 0 & 5 & 5 \\
0 & 0 & 0 & 2 & 2
\end{bmatrix}
=
\begin{bmatrix}
0.14 & 0 \\
0.42 & 0 \\
0.56 & 0 \\
0.70 & 0 \\
0 & 0.59 \\
0 & 0.74 \\
0 & 0.29
\end{bmatrix}
\times
\begin{bmatrix}
12.36 & 0 \\
0 & 9.48
\end{bmatrix}
\times
\begin{bmatrix}
0.57 & 0.57 & 0.57 & 0 & 0 \\
0 & 0 & 0 & 0.70 & 0.70
\end{bmatrix}
$$

Movies: Matrix, Alien, Serenity, Casablanca, Amelie

users

rank = 2

7X5   7X2   2X2   2X7

Users
Movies
Concepts

# SVD Example: Users-to-Movies

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

*Unique Decomposition*

**SciFi-concept**

**Romance-concept**
(note the **multiplication by -1**)

$$
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
3 & 3 & 3 & 0 & 0 \\
4 & 4 & 4 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 4 & 4 \\
0 & 0 & 0 & 5 & 5 \\
0 & 0 & 0 & 2 & 2
\end{bmatrix}
=
\begin{bmatrix}
0.14 & 0 \\
0.42 & 0 \\
0.56 & 0 \\
0.70 & 0 \\
0 & -0.59 \\
0 & -0.74 \\
0 & -0.29
\end{bmatrix}
\times
\begin{bmatrix}
12.36 & 0 \\
0 & 9.48
\end{bmatrix}
\times
\begin{bmatrix}
0.57 & 0.57 & 0.57 & 0 & 0 \\
0 & 0 & 0 & -0.70 & -0.70
\end{bmatrix}
$$

Matrix, Alien, Serenity, Casablanca, Amelie

users

37

# SVD Example: Latent Factors

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



the colors

$$
\begin{array}{c|ccccc}
 & \text{Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\
\hline
 & 1 & 1 & 1 & 0 & 0 \\
 & 3 & 3 & 3 & 0 & 0 \\
 & 4 & 4 & 4 & 0 & 0 \\
\text{users} & 5 & 5 & 5 & 0 & 0 \\
 & 0 & 0 & 0 & 4 & 4 \\
 & 0 & 0 & 0 & 5 & 5 \\
 & 0 & 0 & 0 & 2 & 2 \\
\end{array}
$$

= U $\Sigma$ $V^T$

$D = 3$
$d = 2$

$u_2$, $u_1$

**"Concepts"**
a.k.a. **Latent dimensions**
a.k.a. **Latent factors**

38

# SVD Example: Beyond Blocks

*Assuming we changed A*

*Values of U changed a little*

*Value of $\Sigma$ & $V^T$ changed, a new row introduced*

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

7X5

=

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

7X3

X

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

3X3

X

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

3X5

*Now, rank(A) = 3*

39

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

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

**SciFi-concept**

**Romance-concept**

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

Matrix A — users / movies

**0.13** reflects the amount of similarity bet ⟨ 1st user / 1st Concept

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

SciFi-concept · Romance-concept

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

=

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

* when interpreting values, don't forget reflecting on all the multiplications $U\Sigma V^T$

SVD decomposition with SciFi-concept and movie ratings.

Column labels: Matrix, Alien, Serenity, Casablanca, Amelie

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

SciFi-concept

"strength" of the SciFi-concept

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

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

Matrix  Alien  Serenity  Casablanca  Amelie

SciFi-concept

SciFi-concept

# SVD: Interpretation

*Such decomposition allows for unvealing the true latent dimensionality/rank, allowing for better interpretability of the factors $U$, $\Sigma$ & $V^T$.*

- $A = U\Sigma V^T$

'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

# Roadmap

- 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 (Non-linear Dimensionality Reduction)

# Recap: Dim. Reduction by Low Rank Approx.

- **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

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

GOAL → finding 2 basis vectors (2, because rank=2)
& for each, getting the new coord [slide 28]

# SVD: Alternative Interpretation

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

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

first right singular vector

$v_1$

Movie 2 rating

Movie 1 rating

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

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

$v_1$
$v_2$
$v_3$
New axes

- **U Σ**: **Gives the <u>coordinates of the points</u> in the projection axis**

"new coord"

**Projection of users on the "Sci-Fi" axis U Σ:**


Movie 2 rating / Movie 1 rating — first right singular vector $v_1$

$$
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
3 & 3 & 3 & 0 & 0 \\
4 & 4 & 4 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 2 & 0 & 4 & 4 \\
0 & 0 & 0 & 5 & 5 \\
0 & 1 & 0 & 2 & 2
\end{bmatrix}
\begin{bmatrix}
1.61 & 0.19 & -0.01 \\
5.08 & 0.66 & -0.03 \\
6.82 & 0.85 & -0.05 \\
8.43 & 1.04 & -0.06 \\
1.86 & -5.60 & 0.84 \\
0.86 & -6.93 & -0.87 \\
0.86 & -2.75 & 0.41
\end{bmatrix}
X
\begin{bmatrix}
0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\
0.40 & -0.80 & 0.40 & 0.09 & 0.09
\end{bmatrix}
$$

(8.43) → 5th user enjoys sci-fi

48

# SVD: Best Approximation

- How to <u>find the best approximation</u>?

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

- How to find the best approximation?

✗ Set <u>smallest singular values</u> to <u>zero</u>!

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

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

**NB** Interpretation of 1.3 : how spread data around $v_3$

↳ low variance, hence to be ignored

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

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

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

7X2          2X2

2X7

*  think of A as sum over rank 1 matrices

$$A = \sigma_1 \, u_1 \, v_1^T + \sigma_2 \, u_2 \, v_2^T$$

*by construction, the rank of the new matrix is 2*

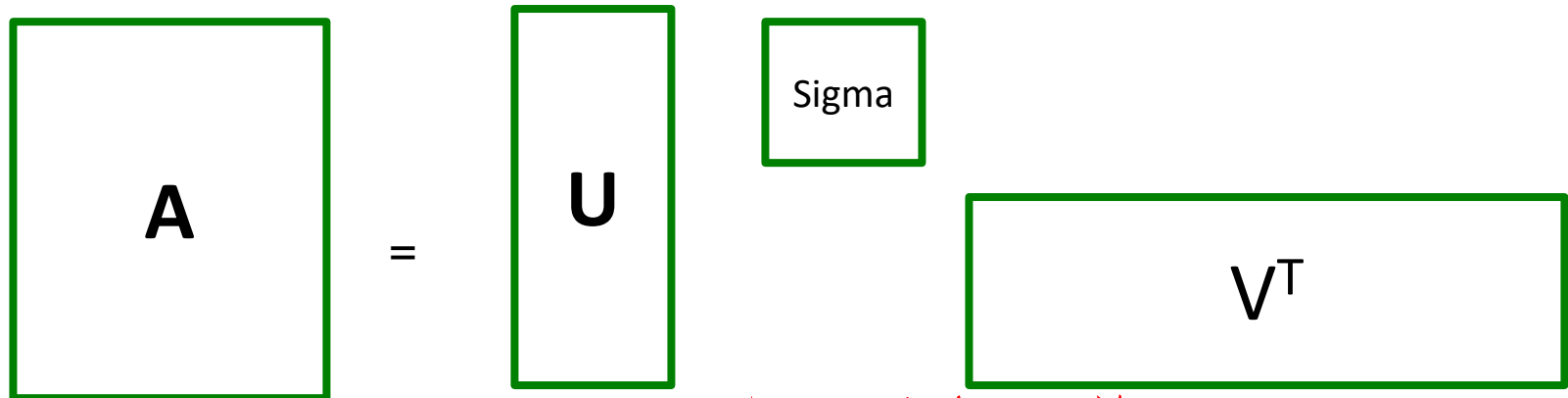$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.92 & 0.95 & 0.92 & 0.01 & 0.01 \\ 2.91 & 3.01 & 2.91 & -0.01 & -0.01 \\ 3.90 & 4.04 & 3.90 & 0.01 &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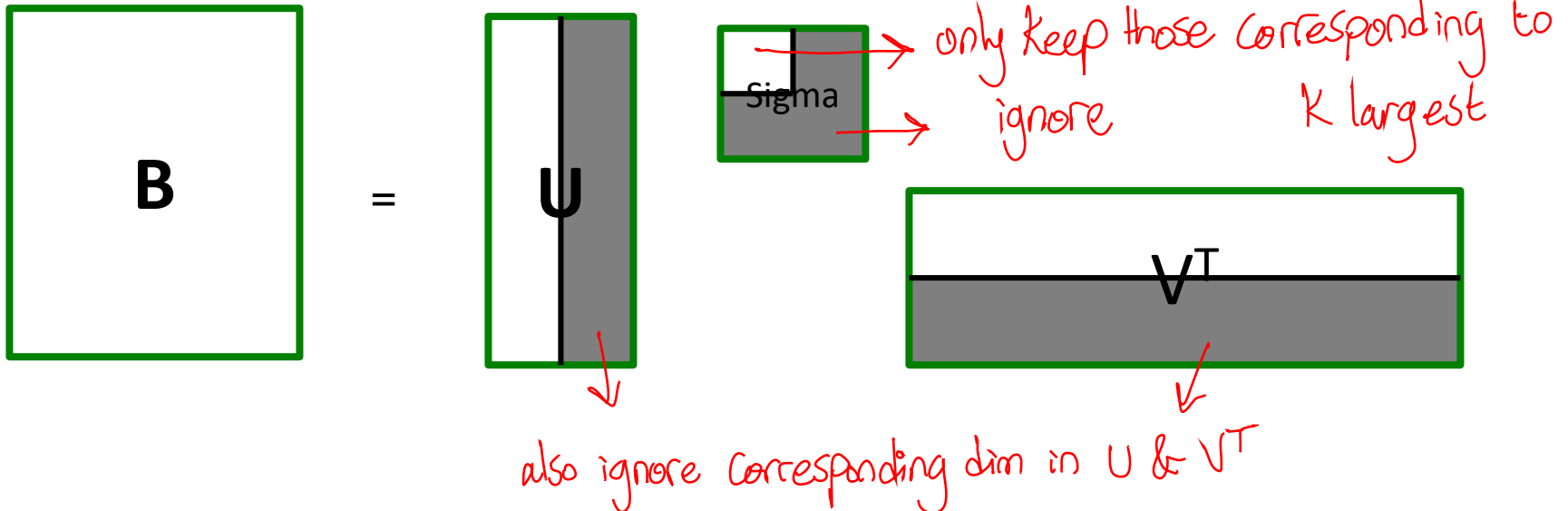
we started with rank 2

the reason for such good approximation is the value we removed is too small
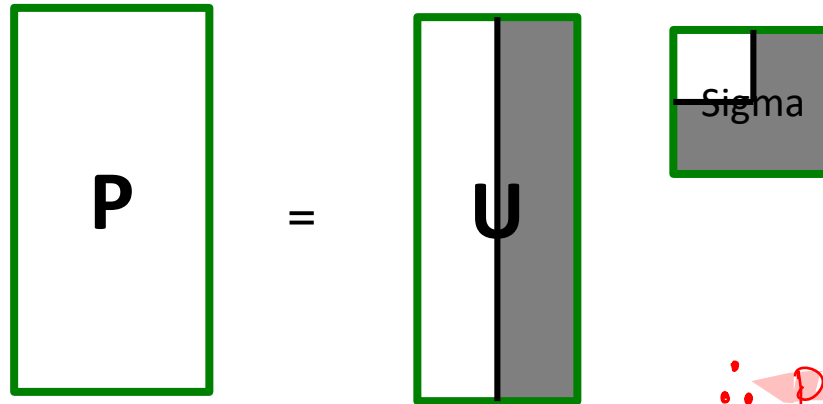
53

# SVD: Best Low Rank Approximation

A = U Sigma V^T

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

$$\text{best rank}(B) = K$$

B = U Sigma V^T

only keep those corresponding to K largest

ignore

also ignore corresponding dim in U & V^T

# SVD: Projection

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



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

$$\therefore P = U\Sigma V^T V = U\Sigma$$

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

Proj data

original data

right singular values

\* No need to calc everything for $U, \Sigma$ & $V^T$, just largest $K$ values

- Recap: Vectors $u_i$ and $v_i$ are <u>of unit length</u>

- <u>W.l.o.g.</u>: $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \ldots \geq 0$

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

✳ Even having larges of $u_i \cdot v_i^T$ would make <u>no</u> diff in terms of error

Having both vectors $u_i$ & $v_i$ on the same scale, the resulting

matrix after multiplication will also be on the same scale.

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \quad \sigma_1 \quad u_1 \quad v^T_1 \; + \quad \sigma_2 \quad u_2 \quad v^T_2 + \ldots$$

$\longleftarrow$ **r terms** $\longrightarrow$

Q: **How many $\sigma_i$ to pick?**
A: Rule of thumb:
keep <u>90% of 'energy'</u>
$\sum_{i=1}^{k} \sigma_i^2 \geq 0.9 \sum_{i=1}^{r} \sigma_i^2$

✳ $\sigma_i$ scales the terms $u_i \cdot v_i^T$

✳ Zeroing small $\sigma_i$ introduces less error

Small $\sigma \longrightarrow$ Small error

57

NB $\sigma_i$ corresponds to the variance

# SVD: Best Low Rank Approximation - Proof

- Theorem: Let A = $U\Sigma V^T$ ($\sigma_1 \geq \sigma_2 \geq \cdots$, rank(A)=r) and B = $USV^T$ with S being a diagonal $r \times r$ matrix where

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

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

  *flatten the vec, calc euclidean dist*

- We have uploaded a detailed proof to the web

  - Note: Many proofs on the web and on other lecture slides are incorrect!

- Remark: B is also an optimal low-rank approximation regarding the spectral norm (operator 2-norm): $\min_B \|A - B\|_2$

  - $\|X\|_2$ = largest singular value of X

# SVD: Best Low Rank Approximation - Proof

- Some facts:

  - $\|X\|_F = \|X^T\|_F$

    - obvious from the definition

  - $\|X\|_F^2 = \text{trace}(X^T X)$                    // trace = sum of diagonal entries

    - easy homework

  - Frobenius norm is invariant to orthonormal transformations $U$

    - Note: If $U^T U = I$ then also $UU^T = I$

    - $\|UX\|_F^2 = \text{trace}((UX)^T(UX)) = \text{trace}(X^T U^T UX) = \text{trace}(X^T X) = \|X\|_F^2$

    - $\|XU\|_F^2 = \|(XU)^T\|_F^2 = \text{trace}(((XU)(XU)^T)$
      $= \text{trace}(XUU^T X^T) = \text{trace}(XX^T) = \|X\|_F^2$

  - Let $A = U\Sigma V^T$ then $\|A\|_F^2 = \|\Sigma\|_F^2 = \sum_i^r \sigma_i^2$

    - follows from above results

# SVD: Complexity

- To compute SVD:

  - $O(n \cdot d^2)$ or $O(n^2 \cdot d)$      (whichever is less)

- But:    #rows      #cols

  - 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 is used

in dimensionality reduction

solving linear sys of equations

$Ax = b$    A being non-square,   underdetermined / overdetermined

# SVD & PCA: Comparison

- Given data X (let's assume it is already centered)
- SVD gives us:
  - X = $U\Sigma V^T$
  - <u>Projected data</u> obtained by $X \cdot V$ (or tr<u>uncated V</u>)
- PCA computes the <u>eigendecomposition</u> <u>of the covariance matrix</u>
  - Covariance matrix: $X^T X$     *assuming X already centered*
  - Eigendecomposition leads to $X^T X = \Gamma \cdot \Lambda \cdot \Gamma^T$
  - P<u>rojected dat</u>a obtained by $X \cdot \Gamma$ (or tr<u>uncated $\Gamma$</u>)
- Let us calculate:
  - $X^T X = \left(U\Sigma V^T\right)^T U\Sigma V^T = V\Sigma^T U^T (U\Sigma V^T) = V\Sigma\Sigma^T V^T = V\Sigma^2 V^T$
  - $\Gamma = V$        **PCA and SVD are equivalent!**
  - $\Sigma^2 = \Lambda$        **squared singular values are variances in new space!**

*V: eigen vec of $X^T X$*

*$\Gamma$: right singular vec of X*

61

# SVD & PCA: Comparison

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

=

find optimal low-rank approximation

(regarding Frobenius norm)    SVD

# Remark: Computation of SVD

We can use the underlined eigendecomposition underlined to calculate underlined the singular value decomposition

- $X^T X = \left(U\Sigma V^T\right)^T U\Sigma V^T = V\Sigma^T U^T (U\Sigma V^T) = V\Sigma\Sigma^T V^T = V\Sigma^2 V^T$

  – V = eigenvectors of $X^T X$

  *right singular vec*

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

  – U = eigenvectors of $XX^T$

  *left singular vec*

- Drawback: Numerically instable

  – better to use specialized algorithms