

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

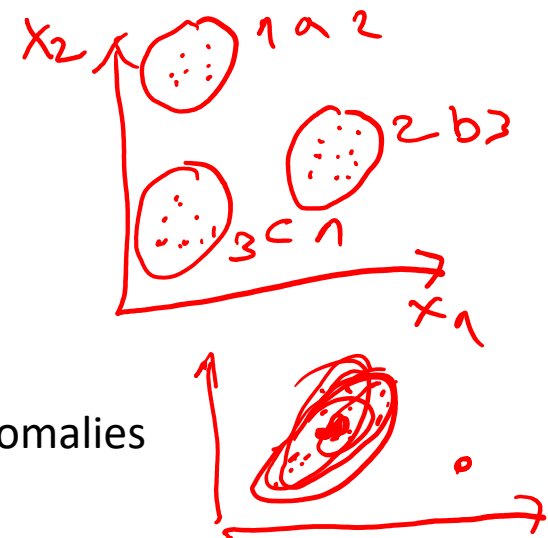
- Chapter: Dimensionality Reduction & Matrix Factorization
 1. **Introduction**
 2. Principal Component Analysis (PCA)
 3. Singular Value Decomposition (SVD)
 4. Matrix Factorization
 5. Neighbor Graph Methods
 6. Autoencoders (Non-linear Dimensionality Reduction)

Introduction: Unsupervised Learning (I)

- Supervised learning aims to map inputs to targets with $y = f(x)$, or in a probabilistic framework it models $p(y|x)$
- Unsupervised learning can be seen as modelling $p(x)$
- We are trying to find the (hidden / latent) structure in the data
 - e.g. find a latent distribution $p(z)$ and a generative transformation $p(x | z)$
we can then 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

$\{x_i, y_i\}$

for each i
 $z_i \sim p(z)$
 $x_i \sim p(x|z_i)$



Introduction: Unsupervised Learning (II)

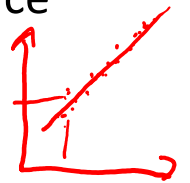
- Unsupervised learning can be viewed as compression
 - compress a data point to a single label corresponding to its cluster
 - compress a data point from a higher dim. to a lower dim. latent space
- Unsupervised learning can be used ...
 - ... as a stand-alone method (e.g. to understand your data, visualization)
 - ... as a 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: $256 \times 256 \times 3 \approx 1.9 \cdot 10^5$
 - 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....)

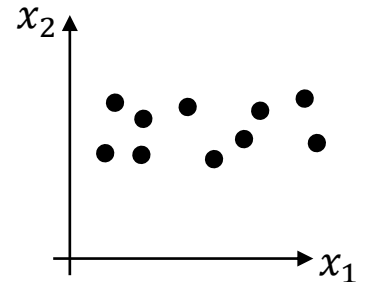


K^2



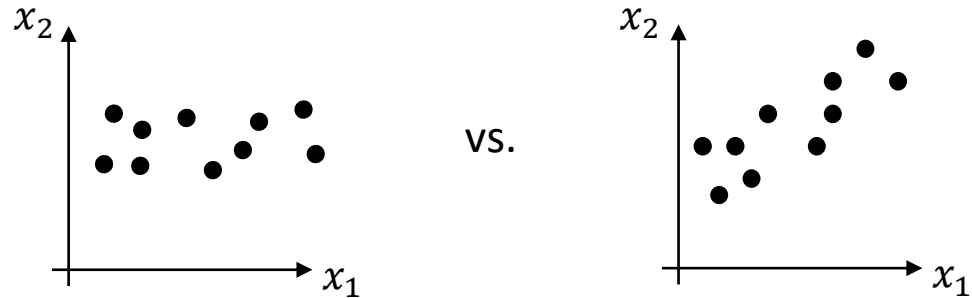
Feature (Sub-)Selection

- 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

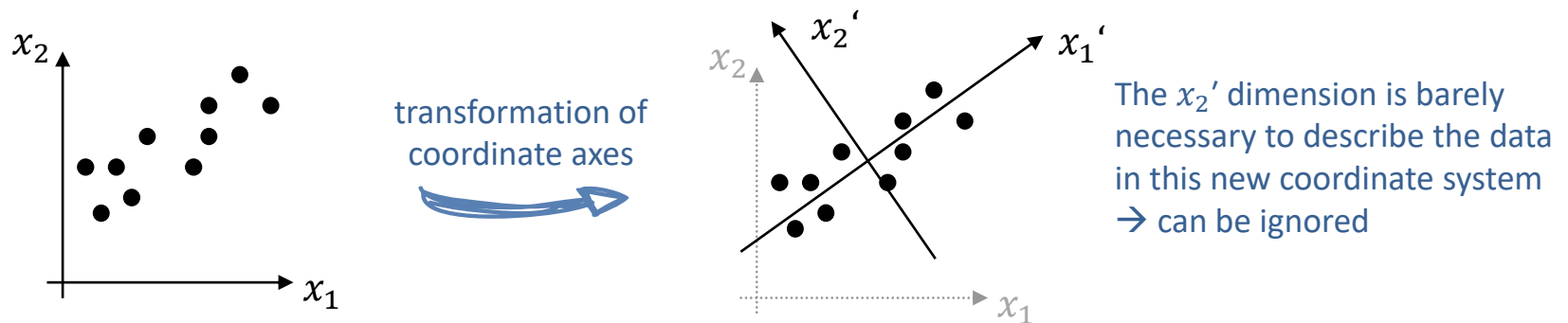


Beyond Feature (Sub-)Selection

- Can we do
 - better?
 - automatic?



- Obviously: Simply discarding whole features not a good idea
 - Features are often correlated



Dim. Reduction via Linear Transformations

- Represent data in a different coordinate system via linear transformations
 - **change of basis** (orthogonal basis transformations)
+ **potentially discarding dimensions**

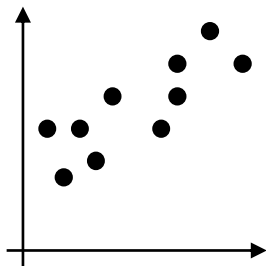
$$F F^T = I$$
$$F^T F = I$$

$$F = \begin{bmatrix} | & | & | \\ \alpha & & \\ | & | & | \end{bmatrix}$$

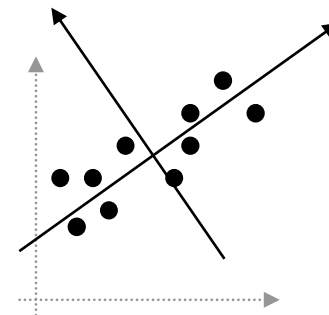
$$\begin{matrix} N \times d \\ \boxed{} \end{matrix} * \begin{matrix} d \times k \\ \boxed{} \\ F \end{matrix} = \begin{matrix} N \times k \\ \boxed{} \\ X' \end{matrix}$$

- Technical:

- use orthonormal transformation matrix $F \in R^{d \times k} < d$
- $(x')^T = x^T \cdot F$ 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



(first center to mean, then
transformation with F)



Discussion: Linear Transformations

- Feature selection is a linear transformation

- What is the matrix F ?

$$F = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}$$

→ 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 $\bar{x}' = \bar{x} \cdot F$
- Let Σ_X be the covariance matrix in the original data space, the covariance matrix in the transformed space is then $\Sigma_{X'} = F^T \cdot \Sigma_X \cdot F$

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)

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



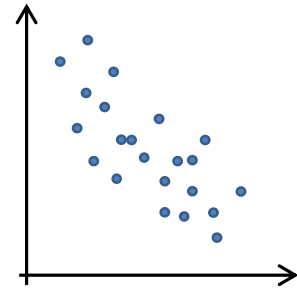
- Approach: Principal Component Analysis (PCA)
 - Find a coordinate system in which the (possibly originally correlated) points are **linearly uncorrelated**
 - The dimensions with no or low variance can then be ignored

Determine the Principal Components

- Goal:

- Transform the data, such that the **covariance between the new dimensions is 0**
- The transformed data points are not linearly correlated any more

$$\text{cov}(\bar{x}_i, \bar{x}_j) = 0 \quad (i \neq j)$$



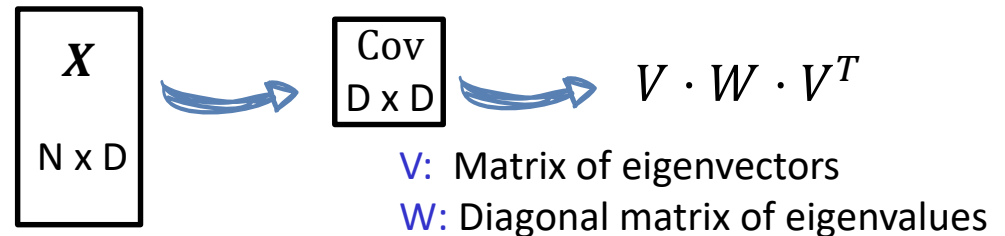
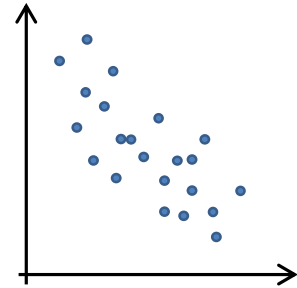
- Given: N d -dimensional data points: $\{\mathbf{x}_i\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{R}^d \forall i \in \{1, \dots, N\}$
- We represent this set of points by a matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$:

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

- 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

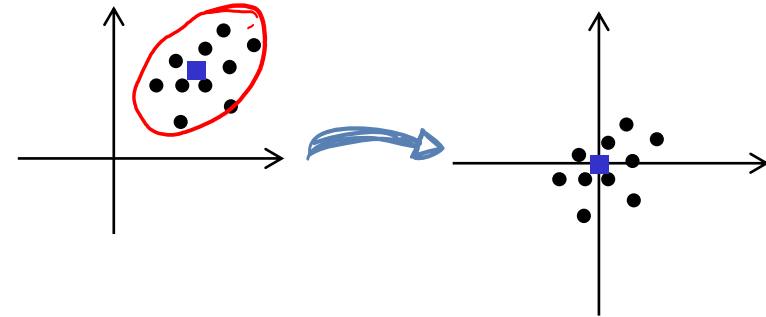
Determine the Principal Components

- 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
 2. Compute the covariance matrix
 3. Use the Eigenvector decomposition to transform the coordinate system



Determine the Principal Components

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



- Shift the points by their mean $\bar{\mathbf{x}} \in \mathbb{R}^d$ (centralized data): $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$

Statistics:

Zero order statistic : number of points N

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

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

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

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

Determine the Principal Components

$$\Sigma_X = \begin{bmatrix} \text{var}(x_{j_1}) & \dots \\ \text{cov}(x_{j_1}, x_{j_2}) & \dots \end{bmatrix}$$

- Determine the variances $\text{Var}(\tilde{X}_j)$ for each dimension $j \in \{1, \dots, d\}$
- Determine the covariance $\text{Cov}(\tilde{X}_{j_1}, \tilde{X}_{j_2})$ between dimensions j_1 and $j_2, \forall j_1 \neq j_2 \in \{1, \dots, d\}$

$$\tilde{X}$$

$$\tilde{X}^t \quad \text{Cov}$$

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

\tilde{X} - data centering

Statistics:

Second order statistic: variance and covariance

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

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

$$X = \begin{bmatrix} | & | \\ \hline & \\ \hline | & | \end{bmatrix} \quad \begin{matrix} j_1 & j_2 \end{matrix}$$

The covariance between dimension j_1 and j_2 is:

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

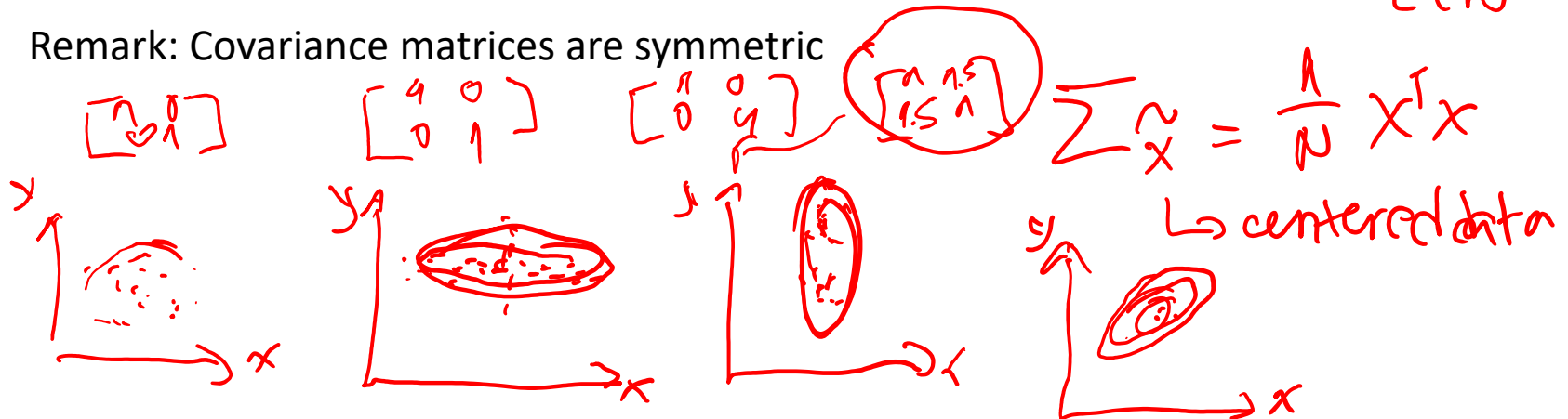
Determine the Principal Components

Statistics (continued):

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

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

- Remark: Covariance matrices are symmetric



Determine the Principal Components



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

- 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: $\Sigma_X = \Gamma \cdot \Lambda \cdot \Gamma^T$

eigenvector of Σ
eigenvalue of Σ

$\Gamma = [\text{columns}]$ $\Lambda = \begin{bmatrix} \circ & & \\ & \circ & \\ & & \circ \end{bmatrix}$

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

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

$$A \gamma = \lambda \gamma$$

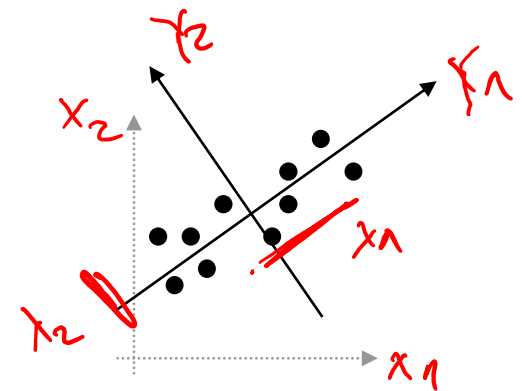
- matrices $\Gamma, \Lambda \in \mathbb{R}^{d \times d}$ with columns of Γ being the normalized eigenvectors γ_i
- Γ is an orthonormal matrix: $\Gamma \cdot \Gamma^T = \Gamma^T \cdot \Gamma = \text{Id}$ ($\Gamma^T = \Gamma^{-1}$)
- Λ is a diagonal matrix with eigenvalues λ_i as the diagonal elements

Determine the Principal Components

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

$$\Sigma_{\tilde{X}} = \underbrace{\Gamma^T}_{\text{side 9}} \tilde{X} \underbrace{\Gamma}_{\text{I}} = \underbrace{\Gamma^T \Gamma}_{\text{I}} \Lambda \underbrace{\Gamma^T \Gamma}_{\text{I}} = \Lambda = \begin{bmatrix} \cdot & & \\ & \ddots & \\ & & \cdot \end{bmatrix}$$

- The **new coordinate system** is defined by the eigenvectors γ_i :
 - Transformed data: $\mathbf{Y} = \tilde{\mathbf{X}} \cdot \underline{\Gamma}$
 - Λ is the covariance matrix in this new coordinate system
 - New system has variance λ_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

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

- Approach

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

- Truncation of Γ

$$\lambda_1 \dots \lambda_k$$

- **Keep only columns (i.e. eigenvectors) of Γ corresponding to the largest k eigenvalues $\lambda_1, \dots, \lambda_k$**

- $Y_{\text{reduced}} = \tilde{X} \cdot \Gamma_{\text{truncated}}$

$$\Gamma = \begin{bmatrix} | & | & | \\ \hline & & \end{bmatrix}$$

k

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_7 & \\ & & \lambda_{0.1} \end{bmatrix}$$

- 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

Complexity



- Complexity of PCA:

$$\underbrace{O(N \cdot d^2)}_{\text{Compute covariance matrix}} + \underbrace{O(d^3)}_{\text{Eigenvalue decomposition}} + \underbrace{O(N \cdot d \cdot k)}_{\text{Project data onto the k-dimensional space}} = \underbrace{O(N \cdot d^2 + d^3)}$$

- Remarks on eigenvalue decomposition:

- Usually we are interested in the reduced data only

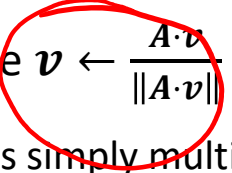

- **Only the k largest eigenvectors required** (i.e. not all of them)

- Use iterative approaches (next slide) for finding eigenvectors

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

- For sparse data even faster: $O(\text{\#it} \cdot \text{\#nz})$ // #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 matrix and v be an arbitrary (normalized) vector
 - Iteratively compute $v \leftarrow \frac{A \cdot v}{\|A \cdot v\|}$ 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

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

How to Compute Eigenvectors?

- Convergence:

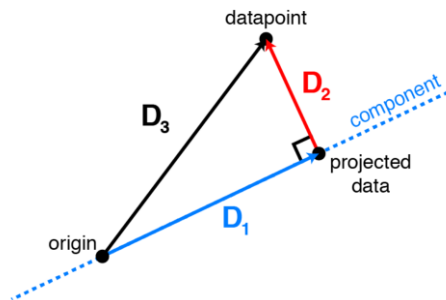
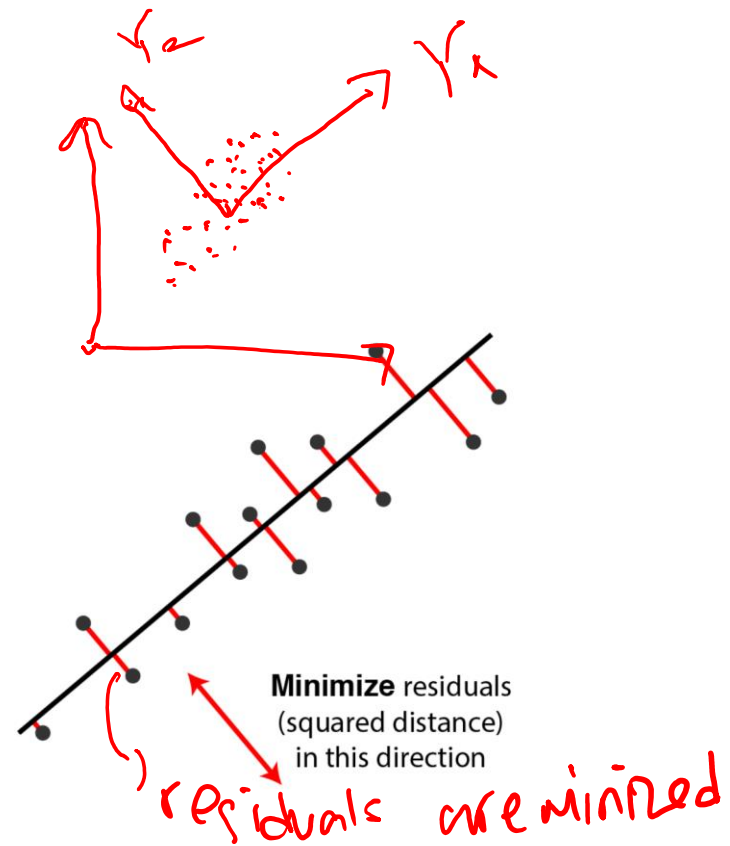
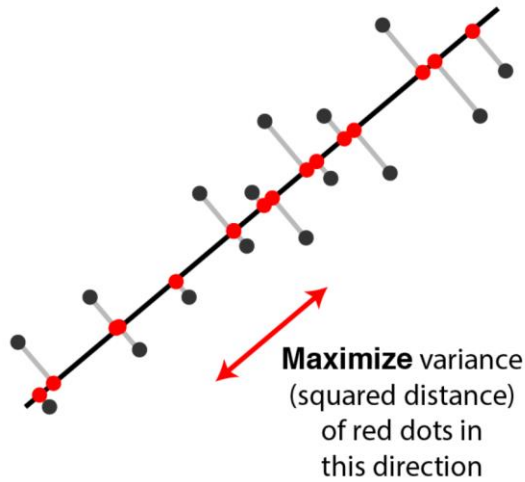
- Linear convergence with rate $|\lambda_2/\lambda_1|$
- Fast convergence if first and second eigenvalue are dissimilar

$O(d^2 \cdot \#it \cdot K)$

- 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: $\hat{A} = A - \lambda_1 \cdot \gamma_1 \cdot \gamma_1^T$
 - \hat{A} has the same eigenvectors as A except the first one has become zero
- Apply power iteration on \hat{A} to find the second largest eigenvector of A

Alternative views of PCA



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

max min

Images adapted from Alexh Williams

PCA vs. Regression

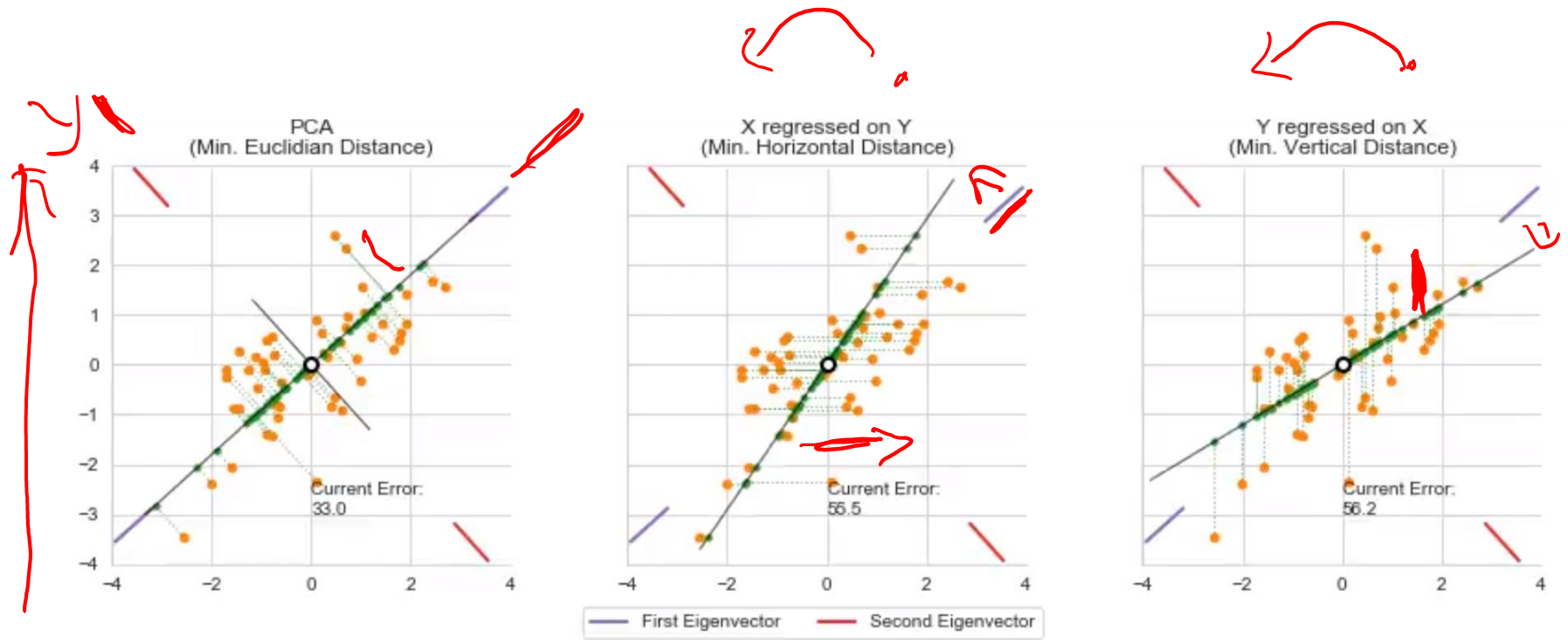
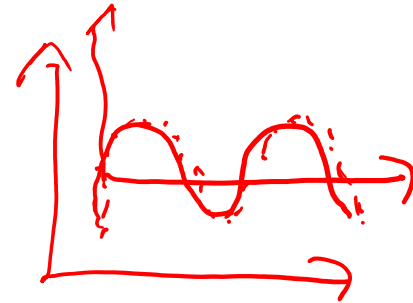


Image adapted from Quentin André

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)

1901

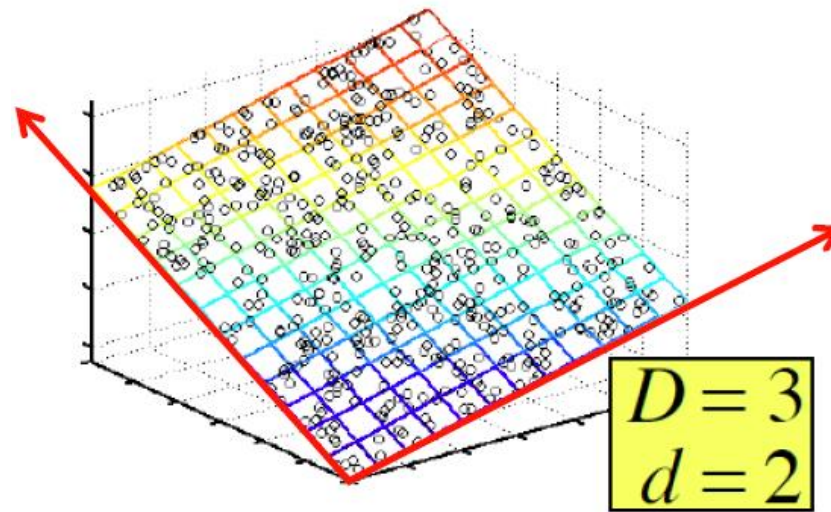


Roadmap

- 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)

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

$$\begin{matrix} a & b & c \\ c = 1 \cdot a + 3b \end{matrix}$$

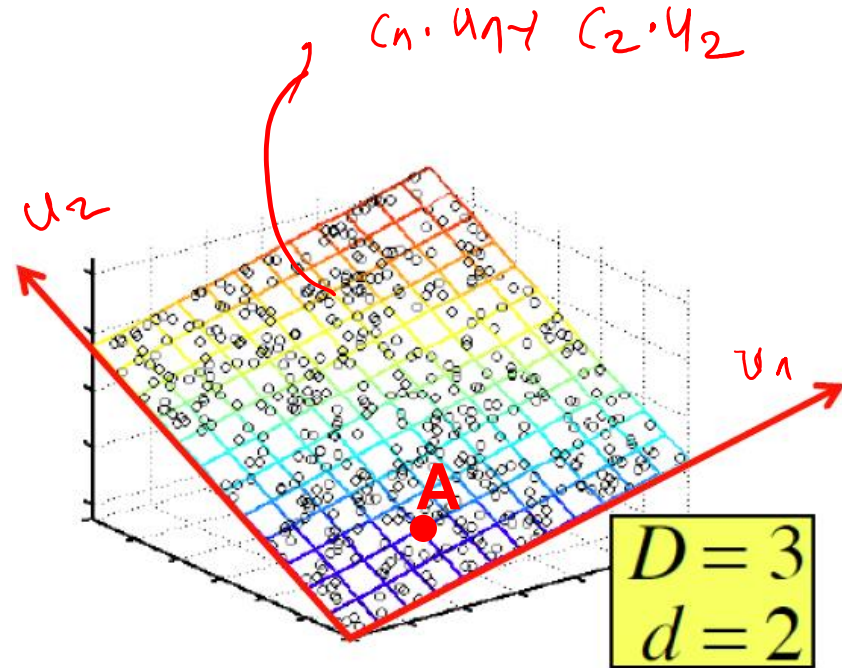
Rank is “Dimensionality”

- Cloud of points in 3D space:

- Think of point positions as a matrix:

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

rank = 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!

$A = 1 \cdot [1 \ 2 \ 1] + 0 \cdot [-2 \ -3 \ 1]$

$0 = 1 \cdot [1 \ 2 \ 1] - 1 \cdot [-2 \ -3 \ 1]$

$= [3 \ 5 \ 0]$

Low Rank Approximation

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

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

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

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

Low Rank Approximation

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

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

- Important: Even though both A and B are $\in \mathbb{R}^{n \times d}$ 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

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

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

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

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

- \mathbf{U}, \mathbf{V} : column orthonormal

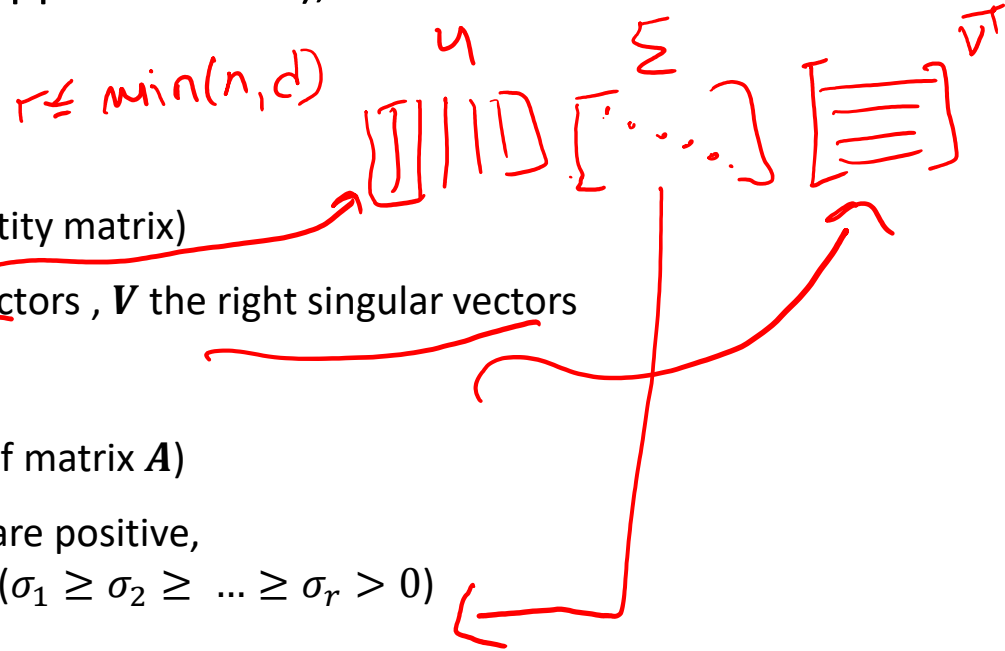
- i.e. $\mathbf{U}^T \mathbf{U} = \mathbf{I}$; $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ (\mathbf{I} : identity matrix)

- \mathbf{U} are called the left singular vectors, \mathbf{V} the right singular vectors

- $\mathbf{\Sigma}$: diagonal

- $r \times r$ diagonal matrix (r : rank of matrix \mathbf{A})

- entries (called singular values) are positive, and sorted in decreasing order ($\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$)



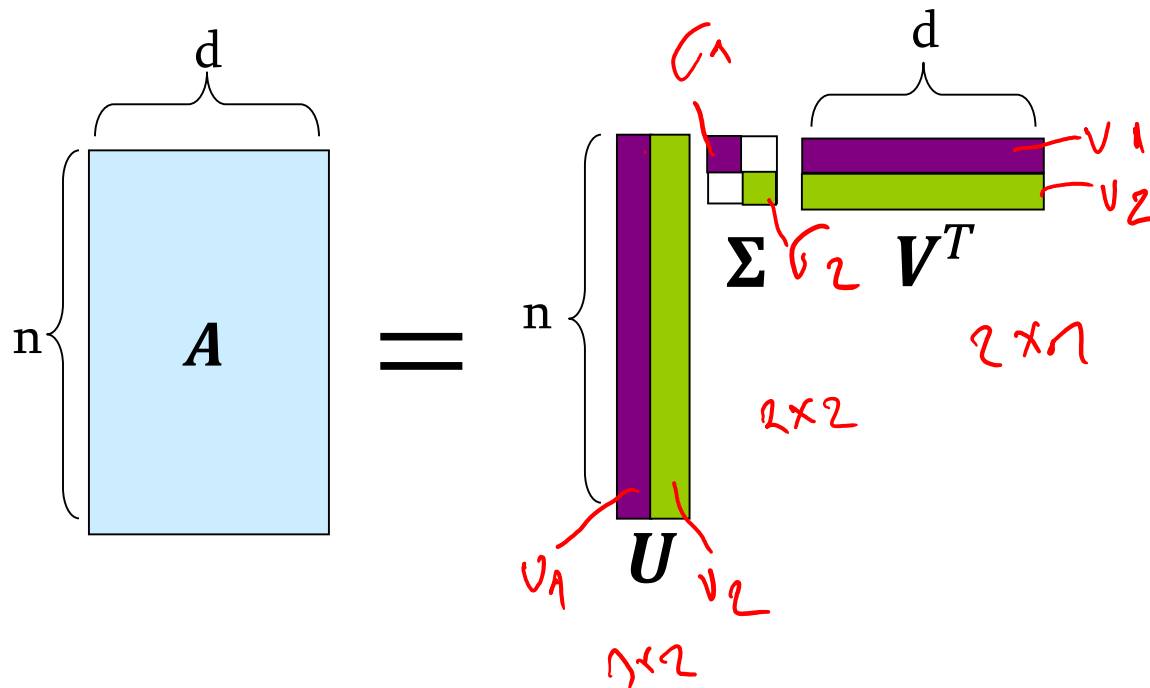
- Remark: The decomposition is (almost) unique

- see e.g. multiplication by -1

Singular Value Decomposition

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

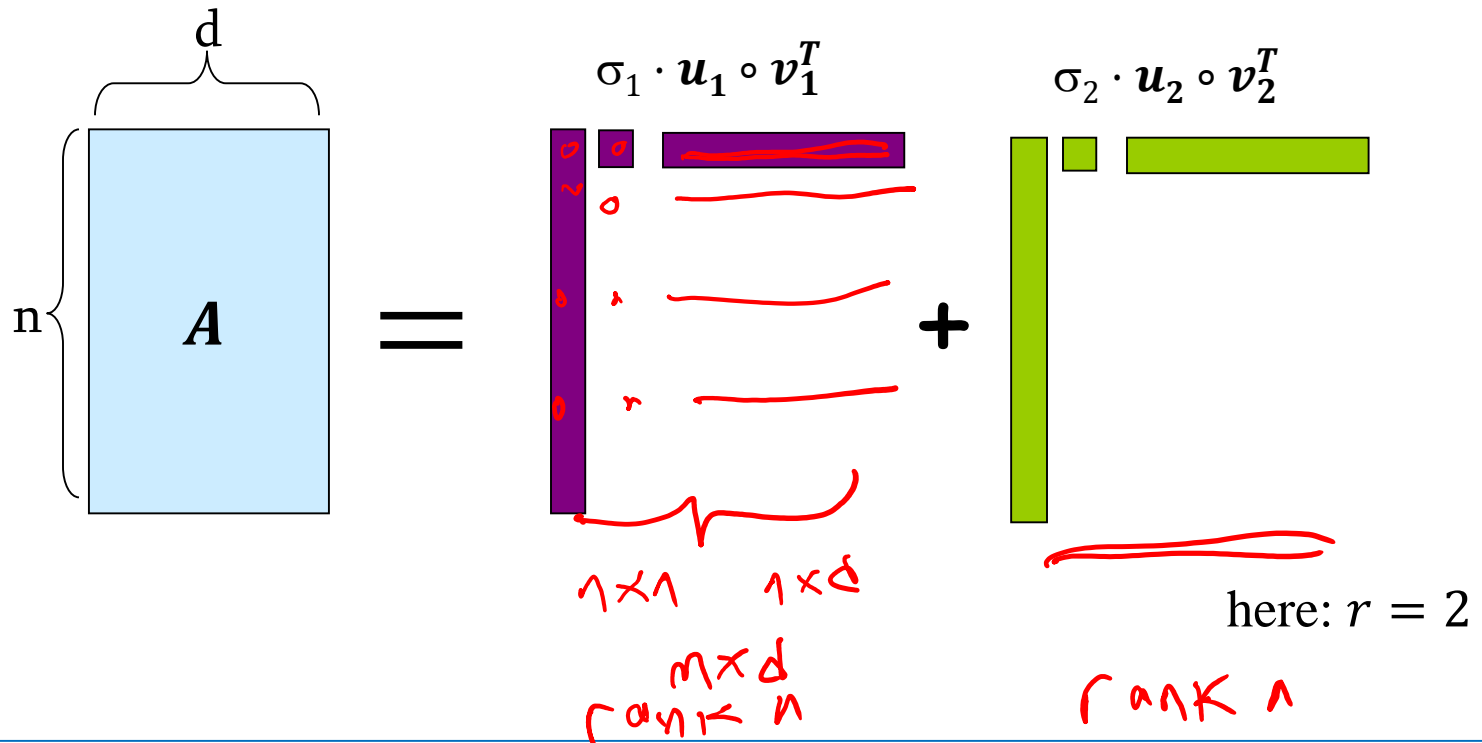
$\text{rank}(A) = 2$



here: $r = 2$

Singular Value Decomposition

$$\underline{\underline{A}} = \underline{\underline{U}} \underline{\underline{\Sigma}} \underline{\underline{V}}^T = \sum_{i=1}^{r=2} \sigma_i \cdot \underline{\underline{u}}_i \circ \underline{\underline{v}}_i^T$$



SVD Example: Users-to-Movies

$$A = \begin{bmatrix} | & | & | & | & | \end{bmatrix} = U \Sigma V^T$$

eigen-faces

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

SciFi-concept Romance-concept

users	Matrix	Alien	Serenity	Casablanca	Amelie
1	1	1	1	0	0
3	3	3	3	0	0
4	4	4	4	0	0
5	5	5	5	0	0
6	0	0	0	4	4
7	0	0	0	5	5
8	0	0	0	2	2

7×5

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

7×2 2×2 2×7

V^T

SVD Example: Users-to-Movies

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

Matrix Alien Serenity Casablanca Amelie

users

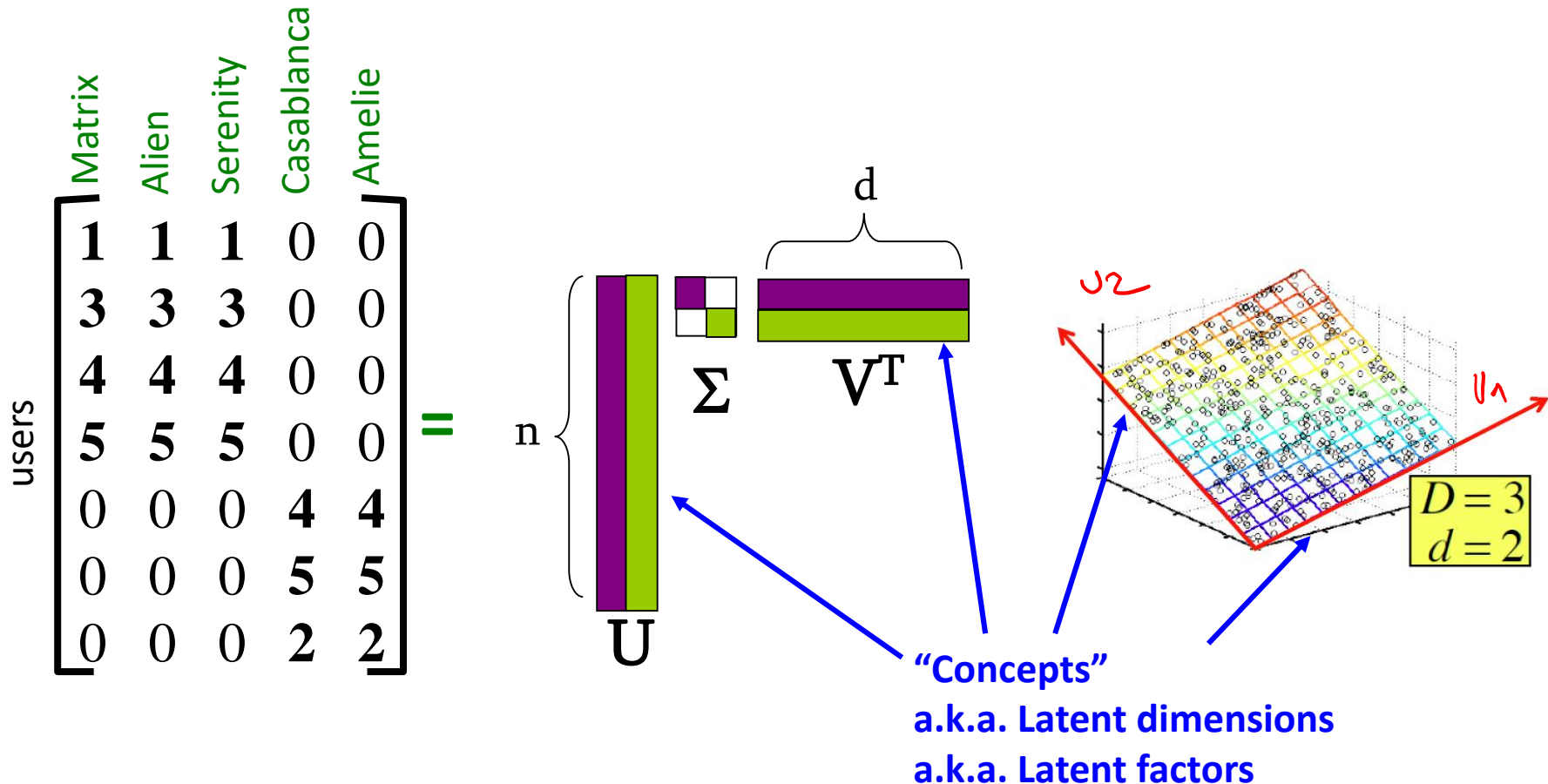
$$\begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 3 & 3 & 3 & 0 & 0 \\
 4 & 4 & 4 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 0 & 0 & 4 & 4 \\
 0 & 0 & 0 & 5 & 5 \\
 0 & 0 & 0 & 2 & 2
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.14 & 0 \\
 0.42 & 0 \\
 0.56 & 0 \\
 0.70 & 0 \\
 0 & -0.59 \\
 0 & -0.74 \\
 0 & -0.29
 \end{bmatrix}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

SciFi-concept

Romance-concept (note the multiplication by -1)

SVD Example: Latent Factors

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



SVD Example: Beyond Blocks

$$\begin{array}{c} \text{Matrix} \\ \text{Alien} \\ \text{Serenity} \\ \text{Casablanca} \\ \text{Amelie} \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & -0.02 & -0.01 \\ 0.41 & -0.07 & -0.03 \\ 0.55 & -0.09 & -0.04 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.07 & 0.73 & -0.67 \\ 0.07 & 0.29 & 0.32 \end{bmatrix} \begin{array}{c} \uparrow \\ \mathbf{X} \end{array} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \begin{array}{c} \mathbf{X} \\ \Sigma \end{array} \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ -0.12 & 0.02 & -0.12 & 0.69 & 0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix} \begin{array}{c} \mathbf{V}^T \end{array}$$

Handwritten annotations: Red circles around the first three rows of the first matrix, the first three columns of the second matrix, and the first three rows of the third matrix. Red arrows point from the first three rows of the first matrix to the first three rows of the second matrix, and from the first three columns of the second matrix to the first three rows of the third matrix. Red 'X' marks are placed below the first three rows of the first matrix and the first three columns of the second matrix. A red 'Σ' is placed below the third matrix.

SVD Example: Beyond Blocks

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

SVD Example: Beyond Blocks

U is “user-to-concept”
similarity matrix

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

SVD Example: Beyond Blocks

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

12.4 9.5 1.3

SVD Example: Beyond Blocks

V is “movie-to-concept”
similarity matrix

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

Diagram illustrating the SVD decomposition of a movie matrix into three components:

- Matrix:** A 7x5 matrix representing movie ratings (rows: Matrix, Alien, Serenity, Casablanca, Amelie, SciFi-concept, SciFi-concept).
- U:** A 7x3 matrix representing movie-to-concept similarity (rows: Matrix, Alien, Serenity, Casablanca, Amelie, SciFi-concept, SciFi-concept).
- S:** A 3x3 diagonal matrix representing singular values.
- V:** A 3x5 matrix representing concept-to-movie similarity (rows: SciFi-concept, SciFi-concept, SciFi-concept).

Annotations:

- Red circles highlight the **SciFi-concept** rows in the U and V matrices.
- Blue arrows point from the **SciFi-concept** rows in the U and V matrices to the corresponding rows in the final product matrix.
- Red circles highlight the **SciFi-concept** rows in the final product matrix.

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
 - Σ : 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

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

$$\begin{aligned} & \begin{bmatrix} 1 & 2 & 1 \end{bmatrix} \\ & \begin{bmatrix} -2 & -3 & 1 \end{bmatrix} \\ & \begin{bmatrix} 1 & 0 \end{bmatrix} \\ & \begin{bmatrix} 0 & 1 \end{bmatrix} \\ & \begin{bmatrix} 1 & -1 \end{bmatrix} \end{aligned}$$

- 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**

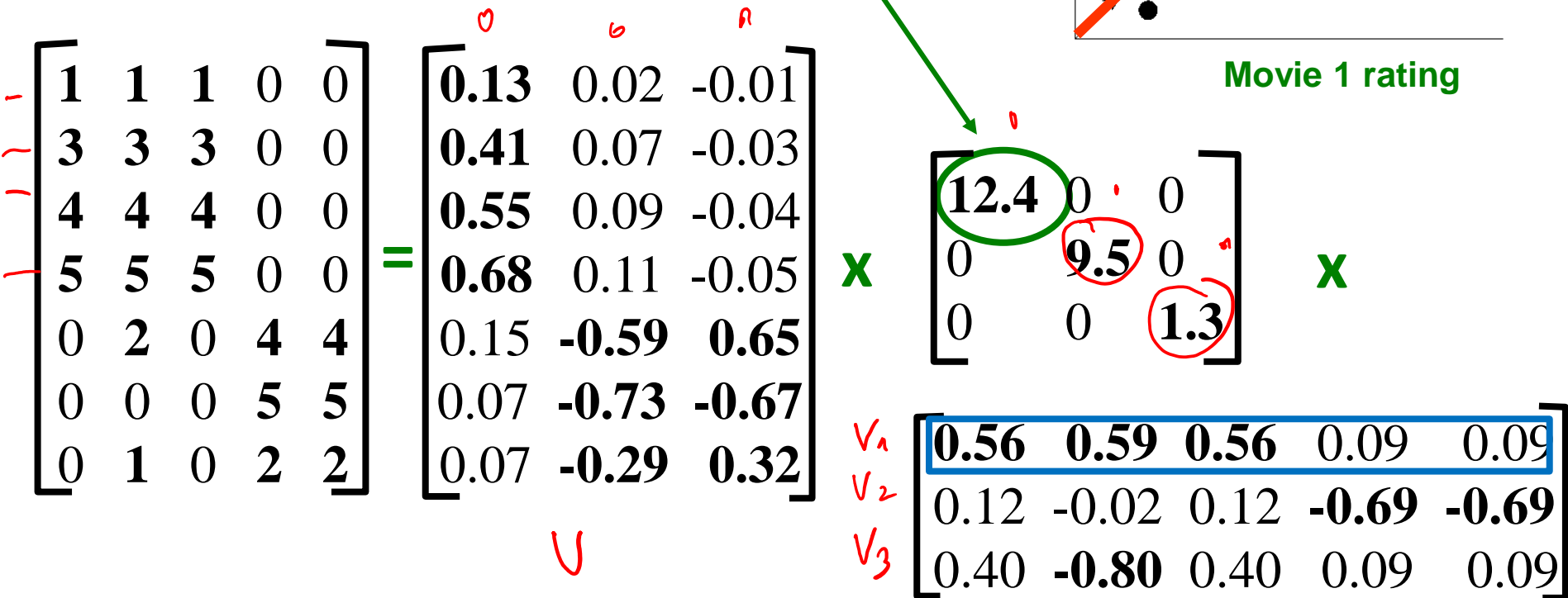
– best = minimize the sum of reconstruction error

– Given matrix $A \in \mathbb{R}^{n \times d}$, find $B \in \mathbb{R}^{n \times d}$ with $\text{rank}(B) = k$ that minimizes

$$\min_{\substack{B \\ \text{rank}(B)=k}} \|A - B\|_F^2 = \sum_{i=1}^N \sum_{j=1}^D (a_{ij} - b_{ij})^2$$

SVD: Alternative Interpretation

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



SVD: Alternative Interpretation

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

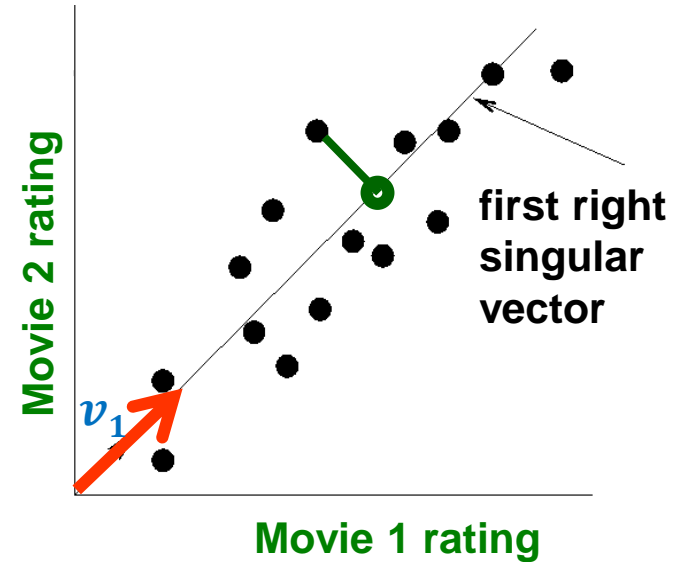
Projection of users on the “Sci-Fi” axis $U\Sigma$:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \times \begin{bmatrix} 1.61 & 0.19 & -0.01 \\ 5.08 & 0.66 & -0.03 \\ 6.82 & 0.85 & -0.05 \\ 8.43 & 1.04 & -0.06 \\ 1.86 & -5.60 & 0.84 \\ 0.86 & -6.93 & -0.87 \\ 0.86 & -2.75 & 0.41 \end{bmatrix}$$

7 x 3

x

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



SVD: Best Approximation

- How to find the best approximation?

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

SVD: Best Approximation

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

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

SVD: Best Approximation

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

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

7x2

x

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

x

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

2x7

SVD: Best Approximation

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

$$A \approx C_1 U_1 V_1^T + C_2 U_2 V_2^T$$

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

A
 U_1
 U_2
 C_1
 C_2
 V_1
 V_2

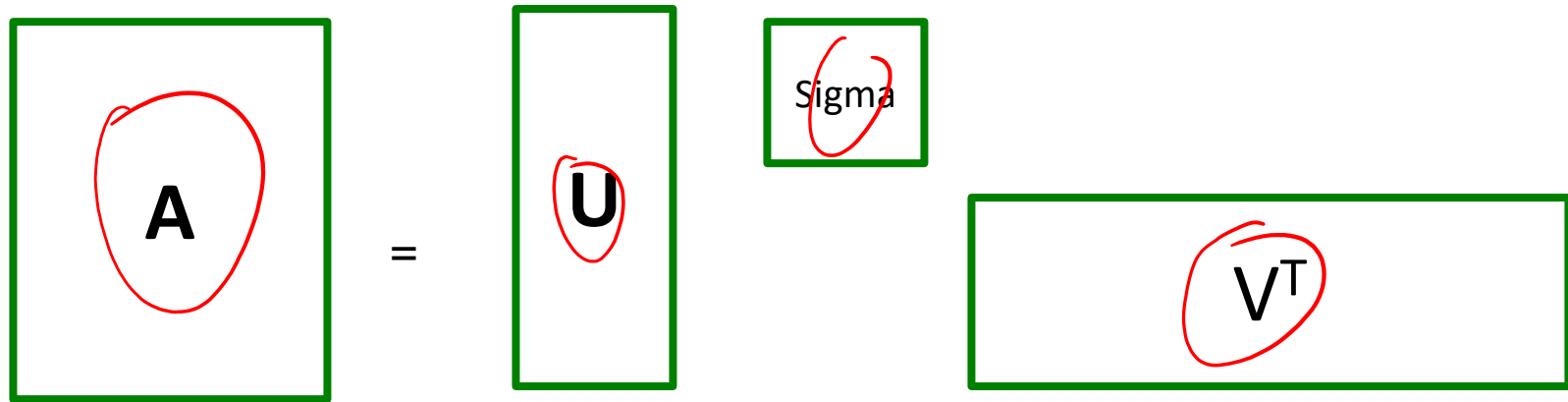
SVD: Best Approximation

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

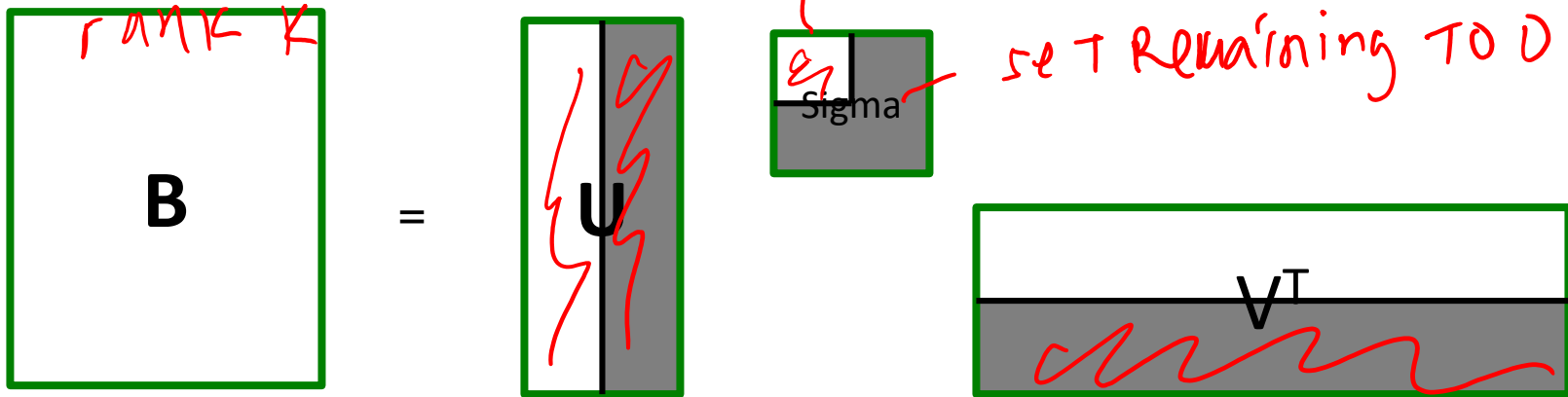
by construction,
the rank of the
new matrix is 2

[illegible]

SVD: Best Low Rank Approximation

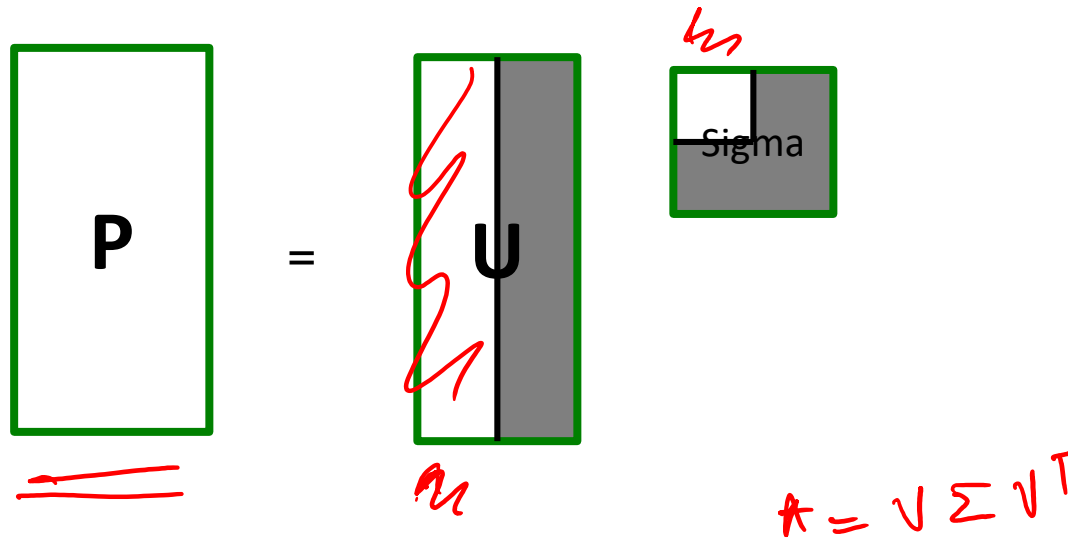


B is best approximation of A



SVD: Projection

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



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

$$P = A V = U \Sigma \underbrace{V^T V}_{I} = U \Sigma$$

Best Approximation – Intuitive Explanation

- Recap: Vectors \mathbf{u}_i and \mathbf{v}_i are of unit length
- W.l.o.g.: $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$

$$\begin{array}{c}
 \begin{bmatrix}
 \mathbf{1} & \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\
 \mathbf{3} & \mathbf{3} & \mathbf{3} & \mathbf{0} & \mathbf{0} \\
 \mathbf{4} & \mathbf{4} & \mathbf{4} & \mathbf{0} & \mathbf{0} \\
 \mathbf{5} & \mathbf{5} & \mathbf{5} & \mathbf{0} & \mathbf{0} \\
 \mathbf{0} & \mathbf{2} & \mathbf{0} & \mathbf{4} & \mathbf{4} \\
 \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{5} & \mathbf{5} \\
 \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{2} & \mathbf{2}
 \end{bmatrix}
 =
 \begin{bmatrix}
 | & | \\
 \mathbf{u}_1 & \mathbf{u}_2 \\
 | & |
 \end{bmatrix}
 \times
 \begin{array}{c}
 \Sigma \\
 \begin{bmatrix}
 \sigma_1 & \text{⊗} \\
 \text{⊗} & \sigma_2
 \end{bmatrix}
 \end{array}
 \times
 \begin{bmatrix}
 \text{---} & \mathbf{v}_1 & \text{---} \\
 \text{---} & \mathbf{v}_2 & \text{---}
 \end{bmatrix}
 \end{array}$$

A
 U
 Σ
 V^T

Best Approximation – Intuitive Explanation

- Recap: Vectors \mathbf{u}_i and \mathbf{v}_i are of unit length
- W.l.o.g.: $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \overset{\text{r terms}}{\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots}$$

σ_i scales the terms $\mathbf{u}_i \cdot \mathbf{v}_i^T$

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

- σ_i scales the terms $\mathbf{u}_i \cdot \mathbf{v}_i^T$
- Zeroing small σ_i introduces less error

SVD: Best Low Rank Approximation - Proof

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

– $s_i = \sigma_i$ for $i=1\dots k$ 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 $\text{rank}(B)=k$

- We have uploaded a detailed proof to the web
 - Note: Many proofs on the web and on other lecture slides are incorrect!
- Remark: 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

$$= \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix}$$
$$= V \Sigma V^T$$

SVD: Best Low Rank Approximation - Proof

- Some facts:

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

$$A = \cancel{U_1} [\cancel{\sigma_1}] + \cancel{U_2} [\cancel{\sigma_2}] - \dots + \cancel{U_k} [\cancel{\sigma_k}] + \dots + U_r [\sigma_r]$$

$$B = \cancel{U_1} [\cancel{\sigma_1}] + \cancel{U_2} [\cancel{\sigma_2}] - \dots + \cancel{U_k} [\cancel{\sigma_k}]$$

$$\|A - B\|_F^2 = \sum_{i=k+1}^r \sigma_i^2 \quad \text{r-k terms}$$

// trace = sum of diagonal entries

$$\Sigma = \begin{bmatrix} \cdot & & \\ & \cdot & \\ & & \cdot \end{bmatrix}$$

$$\|U \Sigma V^T\|_F^2 = \|\Sigma\|_F^2 = \|\Sigma V^T\|_F^2$$

SVD: Complexity

- To compute SVD:
 - $O(\underline{n \cdot d^2})$ or $O(\underline{n^2 \cdot d})$ (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 ...

$$V^+ = V^{-\wedge}$$
$$V^T = V^{-\wedge}$$

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

$$A^\# = V \Sigma^{-\wedge} U^T$$

Moore penrose
pseudo inverse

solve for x

$$\underset{\text{known}}{A} \underset{\text{known}}{x} = \underset{\text{known}}{b}$$
$$x = A^\# b$$

underdetermined

$$\min_x \|x\|_2 \quad Ax = b$$

overdetermined

$$\min_x \|Ax - b\|_2$$

SVD & PCA: Comparison

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

Remark: Computation of SVD

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

$d \times d$

- $X^T X = (U \Sigma V^T)^T U \Sigma V^T = V \Sigma^T U^T (U \Sigma V^T) = V \Sigma \Sigma^T V^T = \underline{V \Sigma^2 V^T}$
– V = eigenvectors of $X^T X$

$n \times n$

- $XX^T = U \Sigma V^T (V \Sigma^T U^T) = U \Sigma \Sigma^T U^T$
– U = eigenvectors of XX^T \leq

- Drawback: Numerically instable
 - better to use specialized algorithms