

*** Applied Machine Learning Fundamentals ***

Principal Component Analysis

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Winter term 2023/2024



Find all slides on [GitHub](#) (DaWe1992/Applied_ML_Fundamentals)

Lecture Overview

| | |
|----------------|---------------------------------|
| Unit I | Machine Learning Introduction |
| Unit II | Mathematical Foundations |
| Unit III | Bayesian Decision Theory |
| Unit IV | Regression |
| Unit V | Classification I |
| Unit VI | Evaluation |
| Unit VII | Classification II |
| Unit VIII | Clustering |
| Unit IX | Dimensionality Reduction |

Agenda for this Unit

① Introduction

② Maximum Variance Formulation

③ PCA Algorithm

④ PCA Applications

⑤ Wrap-Up

Section: Introduction

Why Dimensionality Reduction?
Data Compression
Data Visualization
What is PCA?

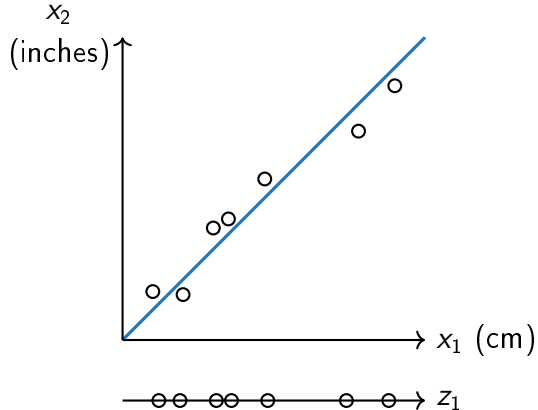
Why Dimensionality Reduction?

- Most data is high-dimensional
- Dimensionality reduction can be used for:
 - **Lossy (!)** data compression
 - Feature extraction
 - Data visualization

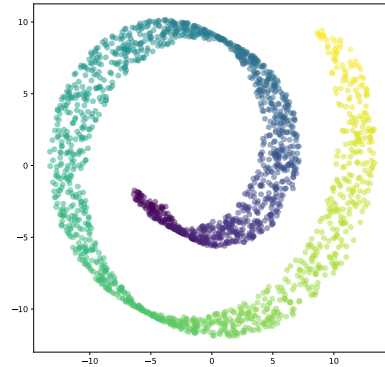
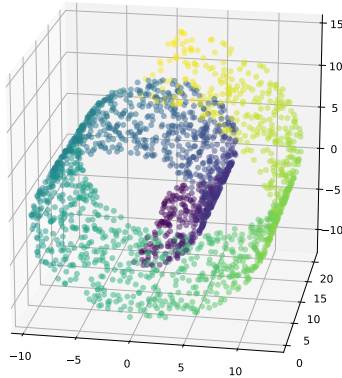
Dimensionality reduction can help to **speed up** learning algorithms substantially. Too many (correlated) features usually **decrease the performance** of the learning algorithm (**curse of dimensionality**).

Use Case I: Data Compression / Feature Extraction

- The features *inches* and *cm* are closely related
- **Problems:**
 - Redundancy
 - More memory needed
 - Algorithms become slow
- **Solution:** Convert x_1 and x_2 into a new feature z_1 ($\mathbb{R}^2 \rightarrow \mathbb{R}$)



Use Case II: Data Visualization



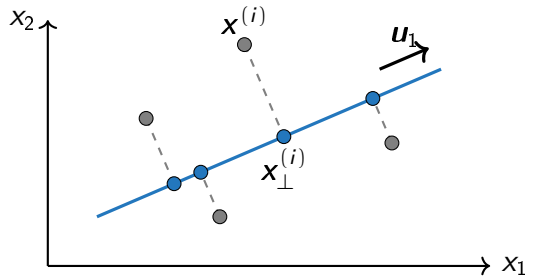
PCA: Principal Component Analysis

- PCA is an **unsupervised** algorithm
- PCA can be defined as the **orthogonal projection** of the data onto a lower dimensional **linear space** (*principal subspace*)
- Consider a dataset of n observations $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$
 - $\mathbf{x}^{(i)} \in \mathbb{R}^m$ ($1 \leq i \leq n$) is an m -dimensional feature vector
 - We want to project the data onto a space having dimensionality $k \ll m$, while **maximizing the variance of the projected data** ($\mathbb{R}^m \rightarrow \mathbb{R}^k$)

Remove dimensions which are the least informative of the data!



Orthogonal Projections (Case: $\mathbb{R}^2 \rightarrow \mathbb{R}$)



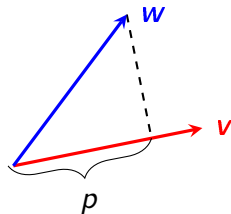
- $\mathbf{x}^{(i)}$ denotes the original data point
- $\mathbf{x}_{\perp}^{(i)}$ is the orthogonal projection of $\mathbf{x}^{(i)}$ onto the vector \mathbf{u}_1

The goal is to find a suitable vector \mathbf{u}_1 so that the variance of the projection is maximized!

Recall: Projection of Vectors

- Let $\mathbf{w}, \mathbf{v} \in \mathbb{R}^2$ be two vectors
- How is the projection of \mathbf{w} onto \mathbf{v} defined?

$$\begin{aligned} p &= \|\mathbf{w}\| \cos \angle(\mathbf{v}, \mathbf{w}) \\ &= \|\mathbf{w}\| \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \cdot \|\mathbf{w}\|} = \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\|} \end{aligned}$$

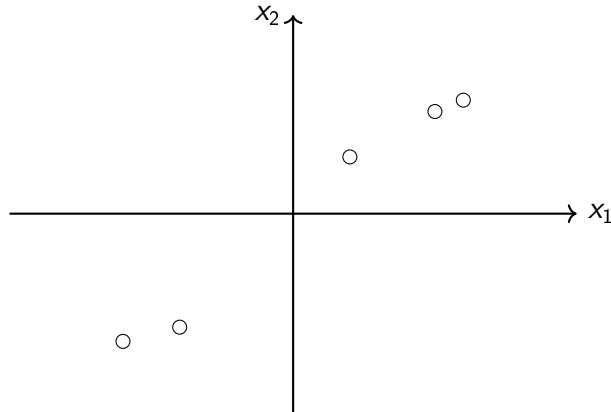


- We will assume \mathbf{u}_1 to be a unit vector, i. e. $\|\mathbf{u}_1\| = 1$
- $\frac{\mathbf{u}_1^T \mathbf{x}^{(i)}}{\|\mathbf{u}_1\|}$ then reduces to the scalar product $\mathbf{u}_1^T \mathbf{x}^{(i)}$

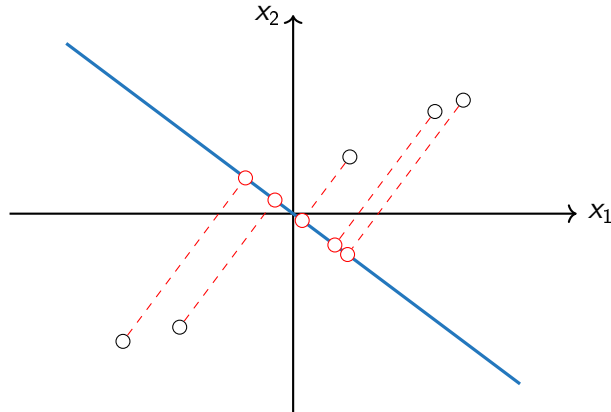
Section: Maximum Variance Formulation

An Example
Formalization of the Problem

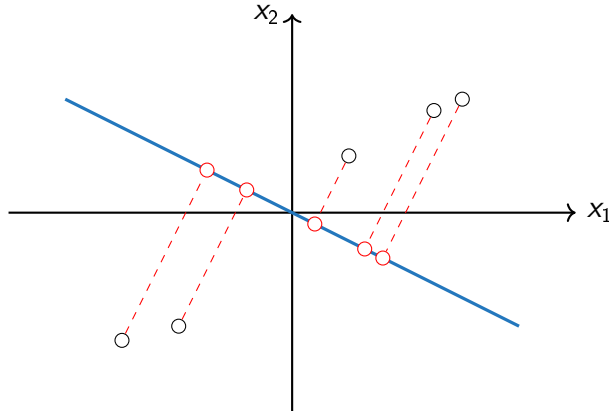
Maximum Variance Formulation



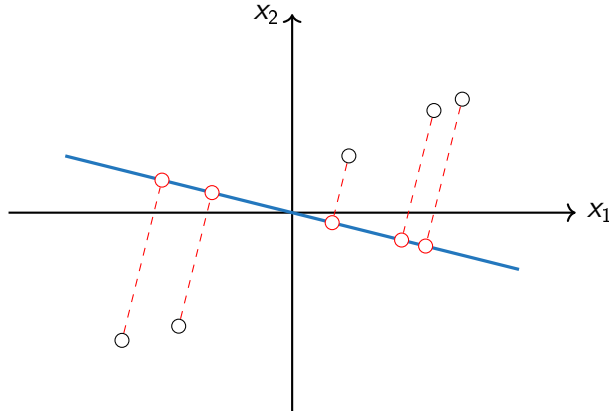
Maximum Variance Formulation



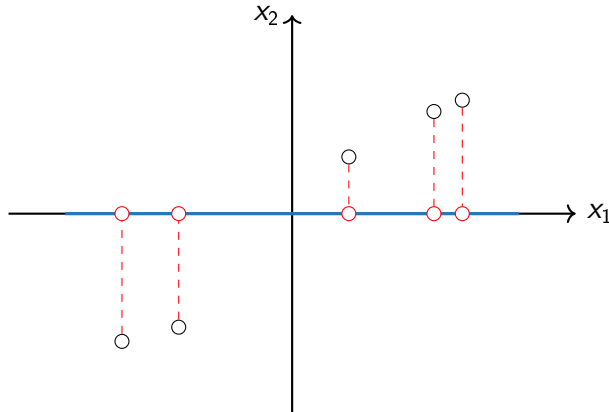
Maximum Variance Formulation



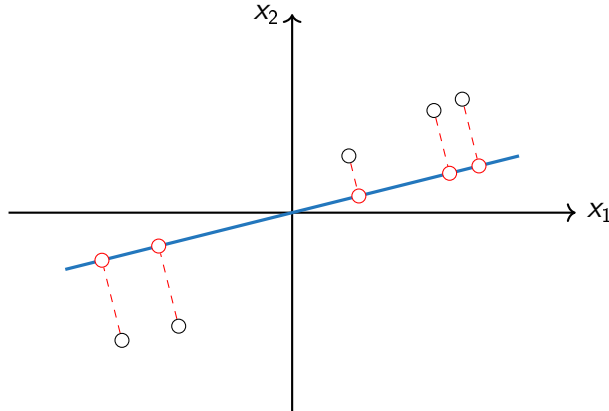
Maximum Variance Formulation



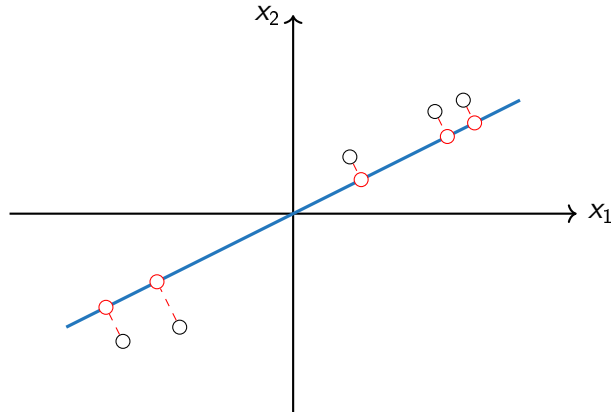
Maximum Variance Formulation



Maximum Variance Formulation



Maximum Variance Formulation





Maximum Variance Formulation (Ctd.)

- In the following we assume $k = 1$
(i. e. we project the data onto a line defined by a unit vector \mathbf{u}_1)
- Each data point $\mathbf{x}^{(i)} \in \mathbb{R}^m$ is projected onto a scalar value $\mathbf{u}_1^\top \mathbf{x}^{(i)} \in \mathbb{R}$
- The mean of the projected data is $\mathbf{u}_1^\top \bar{\mathbf{x}}$, where $\bar{\mathbf{x}} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$
- The variance of the projected data is given by:

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{u}_1^\top \mathbf{x}^{(i)} - \mathbf{u}_1^\top \bar{\mathbf{x}})^2 = \mathbf{u}_1^\top \Sigma \mathbf{u}_1 \quad (1)$$



Maximum Variance Formulation (Ctd.)

- Σ is the covariance matrix defined by:

$$\Sigma := \frac{1}{n} \sum_{i=1}^n \overbrace{(x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^\top}^{\text{Outer product} \rightarrow \text{matrix}} \quad (2)$$

- We have to maximize the projected variance $\mathbf{u}_1^\top \Sigma \mathbf{u}_1$ with respect to \mathbf{u}_1
- **Constraint:** $\|\mathbf{u}_1\| = 1$, otherwise \mathbf{u}_1 grows unboundedly
- We have to solve the following (Lagrangian) optimization problem:

$$\max_{\mathbf{u}_1} \{ \mathbf{u}_1^\top \Sigma \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1) \} \quad (3)$$



Maximum Variance Formulation (Ctd.)

- We have to solve

$$\nabla_{\mathbf{u}_1} \{ \mathbf{u}_1^T \Sigma \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \} \stackrel{!}{=} 0$$

- This leads to the **eigenvalue problem** $\Sigma \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$
- The equation tells us that \mathbf{u}_1 must be an eigenvector of Σ
- If we left-multiply by \mathbf{u}_1^T and use $\mathbf{u}_1^T \mathbf{u}_1 = 1$, we see: $\mathbf{u}_1^T \Sigma \mathbf{u}_1 = \lambda_1$

The variance reaches a maximum, if we set \mathbf{u}_1 equal to the eigenvector having the largest eigenvalue λ_1 . This eigenvector is the first principal component and its eigenvalue λ_1 is the variance retained by it.

Maximum Variance Formulation (Ctd.)

- Additional principal components can be defined in an **incremental fashion**
- Choose each new component such that it **maximizes the remaining projected variance**
- All principal components are **orthogonal to each other**
- Projection onto k dimensions:
 - The lower-dimensional space is defined by the k eigenvectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ of the covariance matrix Σ
 - These correspond to the k largest eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_k$

Section: PCA Algorithm

The Algorithm
An Example
Data Reconstruction
Choice of k



PCA Algorithm

Algorithm 1: PCA Algorithm

Input: Input data $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}) \in \mathbb{R}^{n \times m}$, number of dimensions k

Output: Projected data $\mathbf{Z} \in \mathbb{R}^{n \times k}$

- 1 Compute $\bar{\mathbf{x}} \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$ // sample set mean
 - 2 Compute $\Sigma \leftarrow \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \bar{\mathbf{x}})(\mathbf{x}^{(i)} - \bar{\mathbf{x}})^T$ // covariance matrix
 - 3 Perform singular value decomposition (SVD) for Σ : $(\mathbf{U}, \mathbf{S}, \mathbf{V}) = \text{SVD}(\Sigma)$
 - 4 Select the k eigenvectors with the largest eigenvalues: $\mathbf{U}_k \leftarrow \mathbf{U}_{(:, :k)}$
 - 5 Project the data: $\mathbf{Z} \leftarrow \mathbf{U}_k^T \mathbf{X}$
-

Projection of the Data

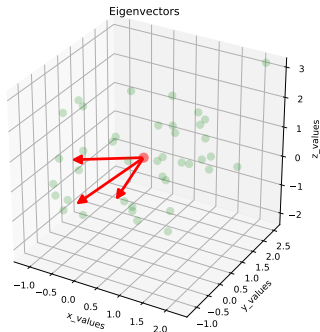
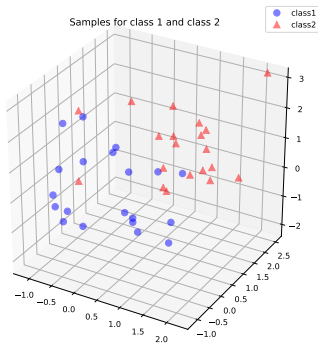
- Matrix \mathbf{U} is obtained by applying **singular value decomposition** to Σ

$$\mathbf{U} = \begin{pmatrix} | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_m \\ | & | & & | \end{pmatrix} \in \mathbb{R}^{m \times m} \quad (4)$$

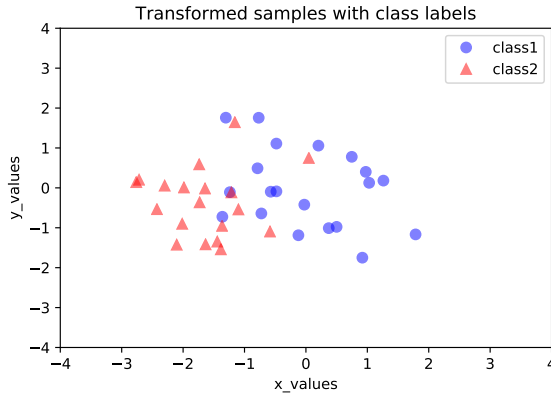
- The projection $\mathbb{R}^m \rightarrow \mathbb{R}^k (k \ll m)$ is performed as follows:

$$\begin{pmatrix} z_1^{(i)} \\ \vdots \\ z_k^{(i)} \end{pmatrix} = \begin{pmatrix} | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_k \\ | & | & & | \end{pmatrix}^\top \begin{pmatrix} x_1^{(i)} \\ \vdots \\ x_m^{(i)} \end{pmatrix} \quad (5)$$

PCA Result



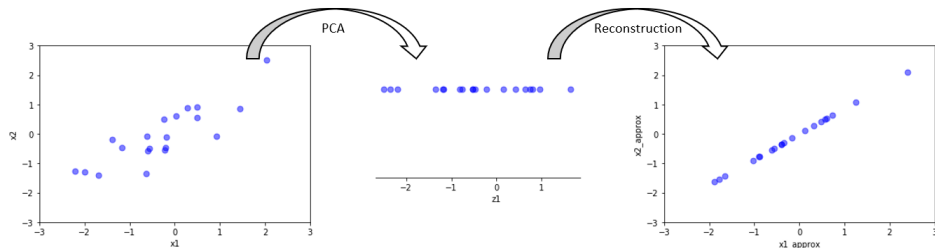
PCA Result (Ctd.)



Reconstruction from compressed Representation

It is possible to compute an **approximate reconstruction** of the data after having applied PCA ($\mathbb{R}^k \rightarrow \mathbb{R}^m$):

$$\mathbf{x}_{\approx}^{(i)} = \mathbf{U}_k \mathbf{z}^{(i)} \quad (6)$$





Choosing the Number of Components

- The goal is to preserve as much variance as possible
- Minimize the **average projection error** given by:

$$\frac{1}{n} \sum_{i=1}^n \| \mathbf{x}^{(i)} - \mathbf{x}_{\approx}^{(i)} \|^2 \quad (7)$$

- **Total variation** in the data is computed as follows:

$$\frac{1}{n} \sum_{i=1}^n \| \mathbf{x}^{(i)} \|^2 \quad (8)$$



Choosing the Number of Components (Ctd.)

- Typically, k is chosen to be the smallest value such that:

$$\frac{\overbrace{\frac{1}{n} \sum_{i=1}^n \| \mathbf{x}^{(i)} - \mathbf{x}_{\approx}^{(i)} \|^2}^{\text{average projection error}}}{\underbrace{\frac{1}{n} \sum_{i=1}^n \| \mathbf{x}^{(i)} \|^2}_{\text{total variation}}} \leq \gamma \quad \gamma \in [0, 1] \quad (9)$$

- This means that $(1 - \gamma) \cdot 100 \%$ of the variance is retained



You can be more efficient...

- The above algorithm is computationally very expensive
- The same result can be computed much more efficient
- Remember: $(\mathbf{U}, \mathbf{S}, \mathbf{V}) = \text{SVD}(\mathbf{\Sigma})$
- We can use the $(m \times m)$ -matrix \mathbf{S}
- \mathbf{S} is a diagonal matrix containing the eigenvalues on its main diagonal:

$$\mathbf{S} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_m \end{pmatrix}$$



You can be more efficient... (Ctd.)

- Sort the eigenvalues in descending order: $(\lambda_1^*, \lambda_2^*, \dots, \lambda_m^*)$
- λ_1^* is the largest eigenvalue, λ_2^* the second-largest, etc.
- Choose the smallest k which satisfies the inequality:

$$\frac{\sum_{j=1}^k \lambda_j^*}{\sum_{j=1}^m \lambda_j^*} \geq \gamma \quad \gamma \in [0, 1] \quad (10)$$

- γ specifies the fraction of variance to be retained overall
- **Advantage:** The matrix \mathbf{S} has to be computed only once and can be reused for all k

Section: PCA Applications

Eigenfaces
Face Morphing

Application of PCA to Images: Eigenfaces



Figure: 100 images of faces



Figure: First 36 principal components

Application of PCA to Images: Eigenfaces (Ctd.)

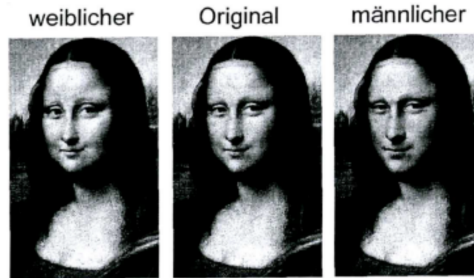


Figure: Original images



Figure: Reconstructed images

Application of PCA to Images: Face Morphing



Section: Wrap-Up

Summary
Self-Test Questions
Lecture Outlook

Summary

- Dimensionality reduction is important to avoid the **curse of dimensionality** 💀...
- ...or simply to **visualize high-dimensional data**
- It is defined as the **orthogonal projection** of the data onto a lower-dimensional (linear) space
- We want to **keep the dimensions with the most variance**
- These dimensions are called **principal components**
- Lots of applications: Eigenfaces, Morphing, ...



Self-Test Questions

- 1 How can PCA be defined?
- 2 What is the geometric relationship between the principal components?
- 3 Outline the PCA algorithm!
- 4 How can you recover the original data? Will you get the exact same data?
- 5 Explain how the number of components / dimensions can be chosen!
- 6 Name some use cases of PCA!

What's next...?



The Exam



Just kidding... (maybe)

Thank you very much for the attention!

Topic: *** Applied Machine Learning Fundamentals *** Principal Component Analysis

Term: Winter term 2023/2024

Contact:

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

daniel.wehner@sap.com

Do you have any questions?