

Business Analytics & Machine Learning

Tutorial sheet 10: Principle Component Analysis

Prof. Dr. Martin Bichler, Prof. Dr. Jalal Etesami
Julius Durmann, Markus Ewert, Johannes Knörr, Yutong Chao

Exercise T10.1 *Principle component analysis*

Given the following dataset $D = \{(-3, -1, -1), (0, -1, 0), (-2, -1, 2), (1, -1, 3)\}$, compute its principal components by following the PCA algorithm introduced in class and generate the transformed data. Each tuple of the set D represents an observation or row vector.

- Calculate the zero-mean dataset X from the given dataset D . Note down the means.
- Calculate the 3×3 covariance matrix Σ_X using the following formulas. What can you infer from it?

$$\text{var}(x_j) = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2$$

$$\text{cov}(x_j, x_k) = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

Reminder: Since the matrix X is centered, you can use the following formulas:

$$\text{var}(x_j) = \frac{1}{N-1} \sum_{i=1}^N x_{ij}^2$$

$$\text{cov}(x_j, x_k) = \frac{1}{N-1} \sum_{i=1}^N x_{ij}x_{ik}$$
- Find the eigenvalues for the covariance matrix by solving the equation: $|\Sigma_X - \lambda I_3| = 0$. You can use the *Laplace expansion* to calculate the determinant.
- Find the corresponding eigenvectors and order them by significance. How is the variance distributed among them?
Hint: Solving the equation $(\Sigma_X - \lambda I_3)v = 0$ gives you the corresponding eigenvectors.
- Compute a one-dimensional PCA projection of the dataset.
- Compute a two-dimensional PCA projection of the dataset.
Hint for e.) and f.): The general formula for projections is: $Z = X\Phi$

Exercise T10.2 *PCA Reconstruction*

Making use of the PCA projections computed in Exercise 10.1, restore the original dataset using the formula: $D \approx Z\Phi^T + \text{means}$

- Reverse the one-dimensional PCA projection to restore the original data. How would the data look when plotted into the original coordinate system?
- What result do you expect when reconstructing the original data from the two-dimensional PCA projection? What is the information loss?

Exercise T10.3 *PCA Understanding*

You are given the following dataset:

$$D = \{(1, -4), (-2, 2), (0, -2), (-1, -1), (2, -3)\}$$

- Apply principal component analysis and determine the two principal component vectors
- Draw the data points, the principal components, and the data set recovered from the 1D projection in a single plot.
- Now another point (a, b) is added to the dataset. As a result, the principal components do not change. Determine a possible point (a, b) .
- We now adjust the original dataset in the following three ways
 - Multiply both coordinates of the first point by some factor $k \in \mathbb{R}$.
 - Multiply all coordinates of all points by the same factor $k \in \mathbb{R}$.
 - Flip all coordinates of all points, e.g. $(1, -4) \rightarrow (-4, 1)$.

In each case, will the eigenvalues and/or principal components change? Discuss your reasoning.

Exercise T10.4 *PCA for image compression*

The goal of this exercise is to use Principal Component Analysis for image compression in Python.

- Load an image of the famous painting "The Starry Night" by Vincent Van Gogh and store red, green, and blue values as color matrices. You can access the painting [here](#). Use the provided template *image_compression_PCA_template.py* to load the painting, store color values, and to plot it.
- Apply a PCA to each color matrix to identify principal color vectors. Choose $n = 5$ as the number of components. Plot the resulting compressed image. The following code snippet may be helpful:

```
for i in range(image.shape[2]):
    # get channel image, normalize data to [0, 1] before applying PCA
    channel_data = img[:, :, i] / 255.0

    # perform pca
    pca = PCA(n_components=number_of_components)
    pca.fit(channel_data)
    compressed_channel_data = pca.inverse_transform(pca.transform(channel_data))
```

- Apply PCA for $n = \{1, 2, 5, 10, 20\}$ (if your computer permits, you can increase n even further). Plot and save the compressed image in each iteration. Observe the size of the image files. Looking at the images, at what point do you clearly identify the painting?
- Determine a reasonable number of clusters using the "elbow criterion". For this purpose, create a scree plot that plots the explained variance of each component against the number of components, e.g. for $n \in [1, 10]$. Does the elbow point correspond to your visual impression in d)?
- Compare the image compression using PCA with the image compression using clustering from last week.