

# Data Mining 2 Homework 1

Jennifer Probst

16 March 2021

## Exercise 1: Proof

Question 1

- We want to show that the objective
 
$$\arg \min_{U \in \mathbb{R}^{d \times r}: U^T U = I} \sum_{i=1}^n \|x_i - U U^T x_i\|_2^2 = \arg \max_{U \in \mathbb{R}^{d \times r}: U^T U = I} \text{trace}(U^T \sum_{i=1}^n x_i x_i^T U)$$
- from lecture slide 19 we find that
 
$$\|x - U U^T x\|_2^2 = \|x\|^2 - x^T U U^T x = \|x\|^2 - \text{trace}(U^T x x^T U)$$
- we find this by  $\|x - U U^T x\|_2^2 = \langle x - U U^T x, x - U U^T x \rangle$ 

$$= \|x\|^2 - 2x^T U U^T x + x^T U U^T U U^T x = \|x\|^2 - x^T U U^T x$$
 now define  $a = U^T x \in \mathbb{R}^{d \times 1} \rightarrow a^T = x^T U \in \mathbb{R}^{1 \times d}$ 
 with hint 1) we can reformulate:  $\|x\|^2 - x^T U U^T x$ 

$$= \|x\|^2 - a^T a = \|x\|^2 - \text{trace}(a a^T) = \|x\|^2 - \text{trace}(U^T x x^T U)$$
 (\*) as  $\|x\|^2$  is constant, we can drop it from the expression to be minimized. independent of  $U$  that is optimized
- if we now sum over components of  $x$ :
 
$$\arg \min_{U \in \mathbb{R}^{d \times r}: U^T U = I} \sum_{i=1}^n \|x_i - U U^T x_i\|_2^2 = \arg \min_{U \in \mathbb{R}^{d \times r}: U^T U = I} \sum_{i=1}^n (\|x_i\|^2 - \text{trace}(U^T x_i x_i^T U))$$
- (\*) 
$$= \arg \min_{U \in \mathbb{R}^{d \times r}: U^T U = I} \sum_{i=1}^n -\text{trace}(U^T x_i x_i^T U) = \arg \max_{U \in \mathbb{R}^{d \times r}: U^T U = I} \sum_{i=1}^n \text{trace}(U^T x_i x_i^T U)$$
 hint 2) 
$$= \arg \max_U \text{trace}\left(\sum_{i=1}^n U^T x_i x_i^T U\right) = \arg \max_U \text{trace}\left(U^T \sum_{i=1}^n x_i x_i^T U\right)$$

## Exercise 2

a) see code

b) see code

c) See the plot of the transformed data in a 2-dimensional subspace of their first two principal components in Figure 1. In the plot the classes are colored by their class label. We observe that the classes are not clearly separated of each other, especially the yellow class is totally mixed into the other ones. All the classes seem to have at least some overlap.

d) The plot of the cumulative variance explained by the principle components can be seen in Figure 3. One can observe that the first few PCs explain almost all the variance. Already 69 percent of the variance (more than 50 percent) is explained by the first PC, the first three principal components explain about 87 percent of the variance and the first five components more than 95 percent of the variance.

e) We perform a data processing step with mean centering of the data and standardization. Looking at the plot of the samples (Figure 2) plotted across their first two components, we find that the data is much better separable, there is only slight overlap between the yellow and green and the green and violet classes. But there are many more PCS needed to explain the variance in the data (Figure 4). We now need 2 PCs to explain 50 percent of the data, three to explain 80 percent of the variance and 17 PCs are needed to explain 95 percent of the variance.

Data normalisation might be a necessary step to properly separate the data. If we don't normalize the data we give features with large values big importance but relative to their large values the variance might not even be so big and the feature not that important.

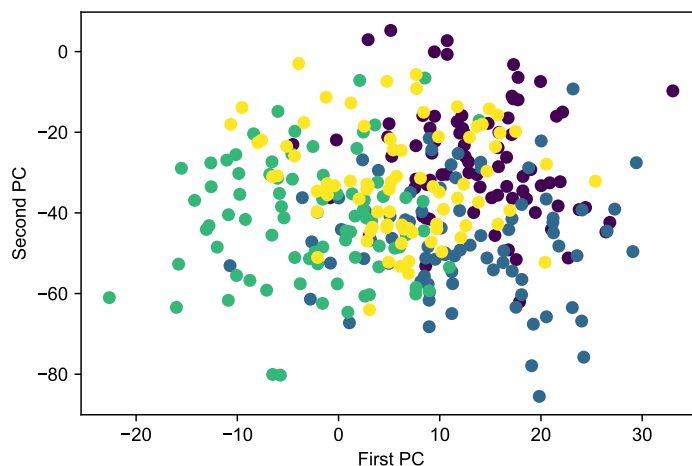


Figure 1: Plot of the transformed data in their first two PCs. Samples are colored according to their class labels.

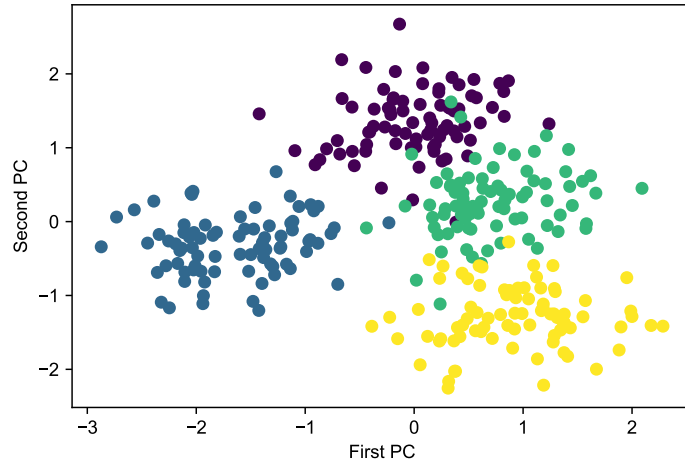


Figure 2: Plot of the normalized and transformed data in their first two PCs. Samples are colored according to their class labels.

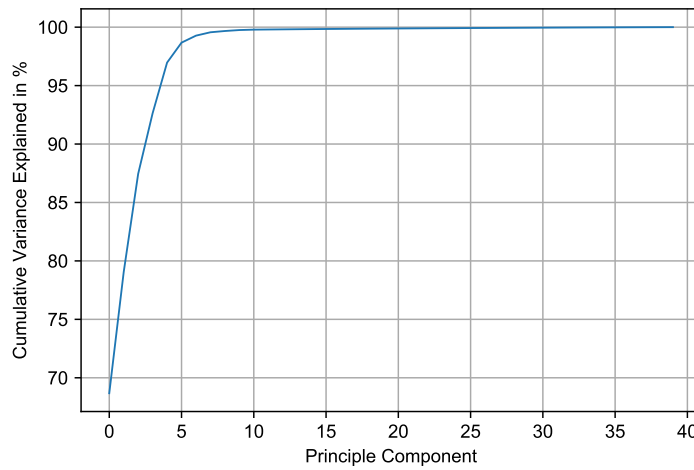


Figure 3: Cumulative variance explained plot for all your principal components for analysis on non-normalized data.

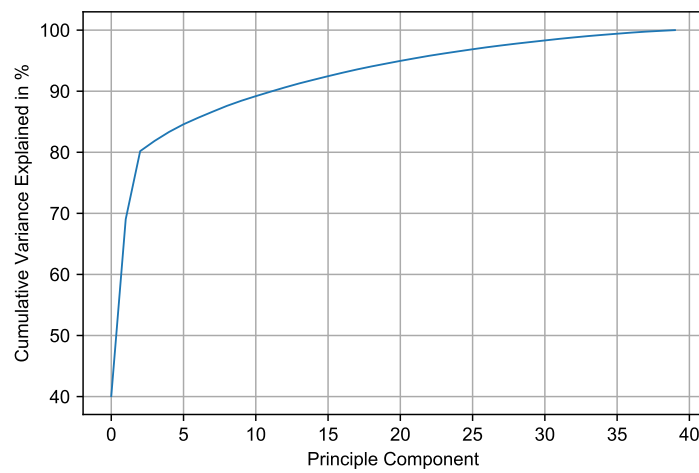


Figure 4: Cumulative variance explained plot for all your principal components for analysis on normalized data.