

sheet03-programming

November 5, 2024

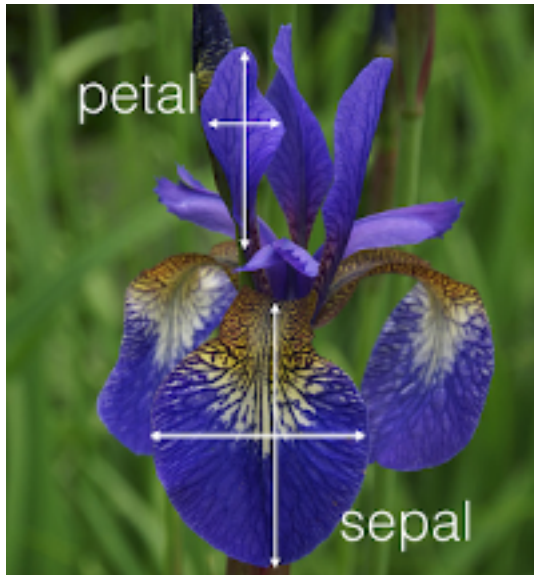
Exercises for the course Machine Learning for Data Science Winter Semester 2024/25

G. Montavon Institute of Computer Science Department of Mathematics and Computer Science
Freie Universität Berlin

Exercise Sheet 3 (programming part)

```
[4]: import numpy
import sklearn
import sklearn.datasets
import matplotlib
%matplotlib inline
from matplotlib import pyplot as plt
```

In this exercise we perform a principal component analysis of the Iris dataset, a famous dataset from Fisher (1936) which one can access from sklearn. Each instance of the dataset is an iris plant which comes with four measurements (1. sepal length in cm, 2. sepal width in cm, 3. petal length in cm, 4. petal width in cm). What these measurements correspond to is depicted in the image below.



In addition to these measurements, the dataset also includes for each instance the type of iris plant (iris setosa, iris versicolour, iris virginica) which we treat here as metadata. Overall, the Iris dataset has 150 instances, and can be stored as an array of size 150 x 4.

The following cell loads the dataset. Additionally, some logarithmic transformation of the input features is performed so that a difference between e.g. 5mm and 6mm is given roughly the same importance as a difference between 5cm and 6cm. Such log-scaling gives more focus on features measuring smaller objects, compared to an approach based on the raw measurements. Also, we add a small increment of 1mm before applying the log-transform to ignore very small quantities that cannot be precisely measured.

```
[5]: dataset = sklearn.datasets.load_iris()

X = numpy.log(0.1+dataset['data'])
T = dataset['target']

target_names = dataset['target_names']
feature_names = dataset['feature_names']

X = X - X.mean(axis=0)

N,d = X.shape
```

Note that one must keep in mind that our features have been log-transformed in any of the subsequent analyses and interpretations. A first basic analysis that can be performed is to measure how much dispersion there is in the data. The total variance of the data can be computed by the code below:

```
[6]: stot = (X**2).mean(axis=0).sum()
print('Total variance: %.3f'%stot)
```

Total variance: 1.027

0.1 Exercise 3 (10 + 10 + 10 P)

We now would like to shed more light into the data by performing a principal component analysis (PCA). The exercise will consist of implementing PCA and then perform various analyses based on the learned model.

0.1.1 Implementing PCA

The simplest way of implementing PCA, is through an eigendecomposition of the data covariance matrix Σ , followed by solving the eigenvalue equation $\Sigma u = \lambda u$.

Task: Implement the PCA algorithm. Your code should return a tuple consisting of (1) a vector L containing the d eigenvalues, in decreasing order, and (2) a matrix U where each column contains the eigenvector associated to the eigenvalue.

Eigenvectors in this matrix should be ordered according to the eigenvalues, i.e. in a way that, for all indices $i = 1 \dots d$, the column $U[:,i]$ contains the eigenvector associated to the eigenvalue $L[i]$.

```
[8]: def PCA(X):
    # Compute the covariance matrix
    cov_matrix = numpy.cov(X, rowvar=False)
```

```

# Eigendecomposition
eigenvalues, eigenvectors = numpy.linalg.eigh(cov_matrix)
# Sort in descending order
sorted_indices = numpy.argsort(eigenvalues)[::-1]
L = eigenvalues[sorted_indices]
U = eigenvectors[:, sorted_indices]
return L, U

# Perform PCA
L, U = PCA(X)

```

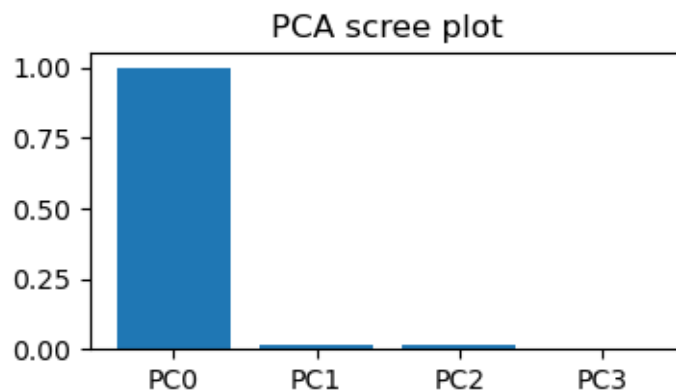
As mentioned in the lecture, PCA can be seen as a way of decomposing the total variance in the data. This decomposition of variance offered by PCA can be rendered in a PCA scree plot. (The following code further verifies that the PCA scree plot corresponds to a decomposition of the total variance, i.e. that the sum of plotted eigenvalues should be equivalent to the measured total variance.)

```

[9]: plt.figure(figsize=(4,2))
plt.title ('PCA scree plot')
plt.bar(numpy.arange(d),L)
plt.xticks(range(d),labels=['PC%d'%i for i in range(d)])
plt.show()

print('Total variance: %.3f'%stot)
print('Sum of eigenvalues: %.3f'%L.sum())

```



```

Total variance: 1.027
Sum of eigenvalues: 1.033

```

We observe that the first principal component captures most of the variance in the data. This suggests that one or a few strongly covariate features capture most of the variations in the data.

0.1.2 Understanding the First Principal Component

We now would like to gain understanding of this highly explanatory first principal component, specifically, to which input features this component responds. For this, we recall that the projection of data on a principal component has the dot-product form

$$z = u^\top x = u_1x_1 + u_2x_2 + \dots + u_dx_d$$

where in our case, u is the principal component stored in our column vector $U[:,0]$.

Task: Write code that displays this projection operation instantiated to our particular problem formulation and PCA model.

```
[10]: # Interpretation of the First Principal Component
print("First Principal Component:")
for i, name in enumerate(feature_names):
    # Each coefficient in U[:, 0] is for the corresponding feature
    print(f"{U[i, 0]:+.3f} * log(0.1 + {name})")
```

```
First Principal Component:
-0.115 * log(0.1 + sepal length (cm))
+0.064 * log(0.1 + sepal width (cm))
-0.561 * log(0.1 + petal length (cm))
-0.817 * log(0.1 + petal width (cm))
```

As can be observed from the formula above, the first principal component responds mostly to a combination of petal length and petal width. It is on the other hand rather insensitive to sepal-related features.

0.1.3 Understanding the First Principal Component

We now would like to get understanding of the data projected on this principal component. The main advantage of looking at the PCA space instead of the feature space is that its dimensionality has been reduced, thereby allowing for using tools such as simple histograms without performing a priori feature selection. The following code projects our data on the first principal component:

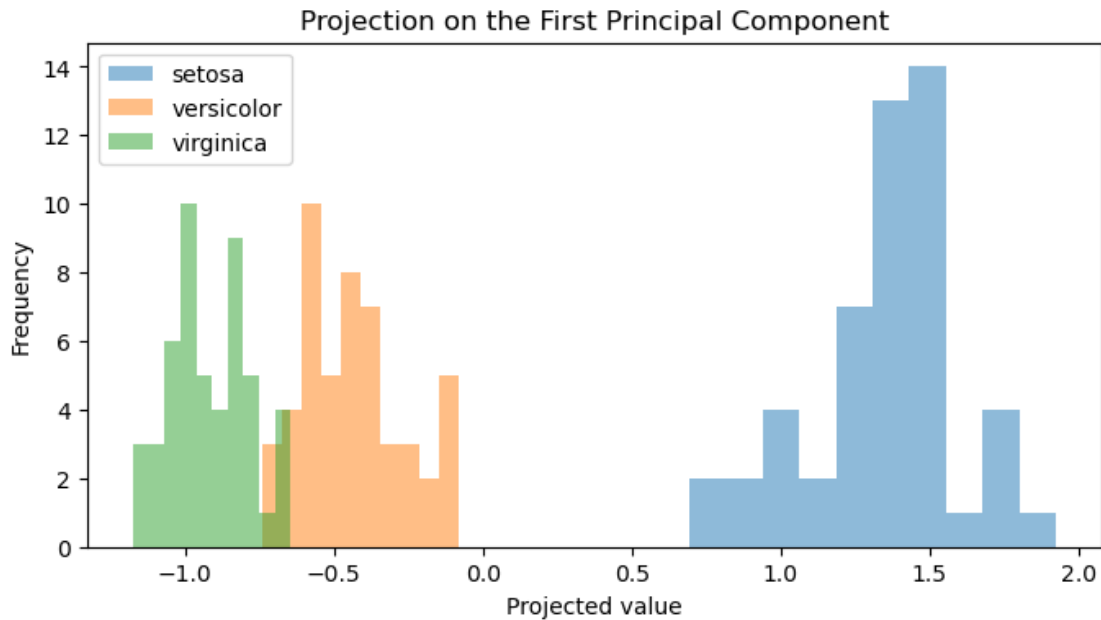
```
[11]: Z = X.dot(U[:,0])
```

Task: Write code that visualizes the projected data in the form of a histogram. Specifically, your visualization should split data according to the meta-data (here, the distinct plant species), and produce for each of them a distinct histogram with a specific color.

```
[12]: # Plot histograms for each species
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 4))
for i, target_name in enumerate(target_names):
    # Create a histogram for each species, using a different color
    plt.hist(Z[T == i], alpha=0.5, label=target_name)
```

```
# Add labels, legend, and title for clarity
plt.legend()
plt.title("Projection on the First Principal Component")
plt.xlabel("Projected value")
plt.ylabel("Frequency")
plt.show()
```



From this analysis, we can observe that the principal component represent meaningful variations in the data, in particular, species appear to be separable. The Iris setosa is separated from the other two species by a large margin, whereas the two other species are also separated to a certain extent, but not fully.