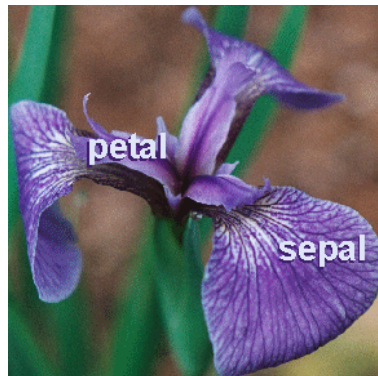


Solutions for exercises marked with a '★' will be made available online, usually in the following week. Only if needed will these exercises be discussed during the tutorial. All other exercises are to be prepared at home and presented by the participants during the tutorial session.

## Exercise 7

### 7.1 Discriminant Analysis

In 1936, Sir Ronald A. Fisher introduced the Iris flower data set as an example of discriminant analysis. It consists of 50 samples of 3 Iris flower species (Iris setosa, Iris virginica and Iris versicolor). 4 features were measured from each sample, being the length and the width of sepal and petal. Based on this, Fisher developed a linear discriminant model to determine which species a sample belongs to.



The data for classes *setosa* and *virginica* are provided by files `setosa.csv` and `virginica.csv`, where each row corresponds to one sample. In order to simplify the problem, we restrict ourselves to the first 2 features only. The goal is to separate both classes by Fisher's Discriminant Analysis (the least-squares version utilizing SVD). Keep in mind, the data must be centered such that the sample mean satisfies

$$\mathbf{m} = \sum_{i=1}^n \mathbf{x}_i = \mathbf{0},$$

where  $\mathbf{x}_i$  is a sample and  $n$  is their total number.

Write a Python script to solve the following problem:

1. Read in both data matrices from files, rows are samples and columns are features.
2. Get  $n$ ,  $n_1$  and  $n_2$ , i.e. the total number of samples and the number of samples per class.
3. Construct the matrix  $A$  containing the samples from both classes. Remember: the columns of matrix  $A$  must be centered.
4. The problem setting is:  $A\mathbf{c} = \mathbf{b}$ .
5. The size of  $\mathbf{b}$  is  $n$ , what is its content?
6. Solve for  $\mathbf{c}$ , do some appropriate plotting.

## 7.2 PCA from scratch\*

For the `hald.csv` data set from last weeks' exercise sheet we now want to program our own implementation of PCA from scratch, i.e. without using Scikit-learn.

Write a Python script to solve the following problem:

- Calculate the empirical covariance matrix.
- Use the eigenvalue decomposition (EVD) to obtain the rotation matrix  $\mathbf{U}$ .
- Is it important whether you perform the eigendecomposition of the covariance matrix before or after standardizing the data?
- Generate appropriate visualizations in order to check whether your implementation performs correctly.

## 7.3 Secondary Structure Prediction\*

The problem of predicting helix structures in protein sequences was introduced in the lecture where we derived an algorithmic variant of classification by least squares. Here, the advantage was that only a “small” matrix  $A^T A$ , vector  $A^T \mathbf{b}$  and scalar  $\mathbf{b}^T \mathbf{b}$  were needed to solve large systems of the form  $A\mathbf{x} = \mathbf{b}$  for  $\mathbf{x}$ .

As the system is overdetermined, our goal is to find the solution vector  $\mathbf{x}$  which minimizes  $\|\mathbf{r}\|^2 = \|A\mathbf{x} - \mathbf{b}\|^2$ . We will analyze how  $\|\mathbf{r}\|^2$  and  $\mathbf{x}^T \mathbf{x}$  change when zeroing a varying number of singular values.

*Please make sure to read the lecture slides!* Given are the two files `AtA.csv` and `Atb.csv` containing the result of the calculations  $\mathbf{A}^t \mathbf{A}$  and  $\mathbf{A}^t \mathbf{b}$  and the additional information  $\mathbf{b}^t \mathbf{b} = 48$ . Write a Python script to solve the following problem:

- For solving  $A\mathbf{x} = \mathbf{b}$  without using  $\mathbf{A}$  or  $\mathbf{b}$  but using  $\mathbf{A}^t \mathbf{A}$  and  $\mathbf{A}^t \mathbf{b}$  instead, we need the singular values  $\mathbf{S}$  and the matrix  $\mathbf{V}$  (for computing  $\mathbf{x} = \mathbf{V}\mathbf{z}$ ). See lecture slides page 155.
- Find  $\mathbf{z}$  which minimizes  $\|\mathbf{S}\mathbf{z} - \mathbf{c}\|^2$
- Compute  $\mathbf{x}$  that corresponds to  $\mathbf{z}$ , i.e.  $\mathbf{x} = \mathbf{V}\mathbf{z}$ . Calculate  $\mathbf{x}^t \mathbf{x}$ .
- Calculate the error for  $\mathbf{x}$  (using  $\mathbf{A}^t \mathbf{A}$ ,  $\mathbf{A}^t \mathbf{b}$  and  $\mathbf{b}^t \mathbf{b}$ ), i.e. calculate  $\mathbf{r}^t \mathbf{r} = \mathbf{x}^t \mathbf{A}^t \mathbf{A} \mathbf{x} - 2\mathbf{x}^t \mathbf{A}^t \mathbf{b} + \mathbf{b}^t \mathbf{b}$ .
- In two different plots, plot  $\mathbf{r}^t \mathbf{r}$  and  $\mathbf{x}^t \mathbf{x}$  as a function of the number of singular values used. There is a total of 21 non-zero values in the matrix  $\mathbf{S}$ .