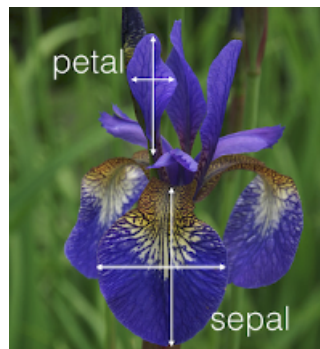


Exercise Sheet 4 (programming part)

```
In [1]: 1 import numpy
        2 import scipy
        3 import sklearn
        4 import sklearn.datasets
        5 import sklearn.decomposition
        6 import torchvision
        7 import matplotlib
        8 %matplotlib inline
        9 from matplotlib import pyplot as plt
```

Exercise 3 (10 + 10 P)

In this exercise we perform a clustering 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 and performs some standardization.

```
In [2]: 1 dataset = sklearn.datasets.load_iris()
        2
        3 X = dataset['data']
        4
        5 X = X - X.mean(axis=0)
        6 X = X / X.std()
```

We focus on one of the simplest clustering methods, which is to find the *connected components* of a graph associated to the dataset. We consider a graph where two nodes are connected if the distance between the corresponding instances (after dataset normalization) is below some threshold value δ .

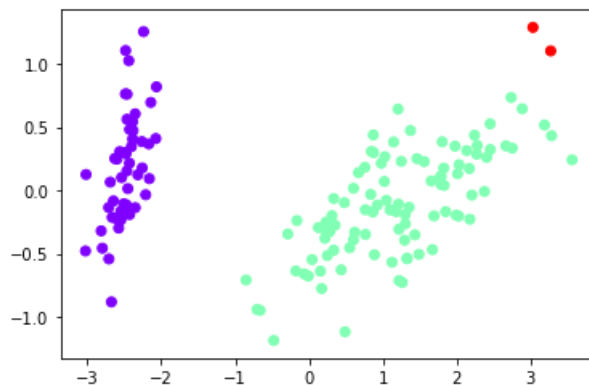
(a) Implement this simple clustering method. (Hint: you can make use of the method `scipy.sparse.csgraph.connected_components` to find the connected components associated to a particular adjacency matrix.)

```
In [3]: 1 class ConnectedComponentsClustering:
2
3     def __init__(self,delta):
4         self.delta = delta
5
6     def fit_predict(self,X):
7
8         # -----
9         # TODO: Replace by your code
10        # -----
11        import solution
12        n_components, Y = solution.cc_fit_predict(self,X)
13        # -----
14
15        return n_components, Y
```

The algorithm is now applied to the Iris dataset with a particular threshold value $\text{delta} = 0.725$. Results are visualized in a PCA plot, where instances are color-coded according to their cluster.

```
In [4]: 1 delta = 0.725
2
3 n_components, Y = ConnectedComponentsClustering(delta).fit_predict(X)
4
5 Z = sklearn.decomposition.PCA(2).fit_transform(X)
6 plt.scatter(*Z.T,c=Y,cmap='rainbow')
```

Out[4]: <matplotlib.collections.PathCollection at 0x7f0a7449d340>



We observe three clusters. Two clusters contain the bulk of the data, and another cluster is very sparsely populated and appears to contain two outliers. Clusters are consistent with what a human can see in the PCA plot.

(b) We now would like to compare the taxonomy derived from the clustering above with the actual iris types available as meta-data. Write code that produces a correspondence table linking the clustering to the classification into iris types.

```
In [5]: 1 T = dataset['target']
2 names = dataset['target_names']
3
4 # -----
5 # TODO: Replace by your code
6 # -----
7 import solution
8 solution.correspondence(names,T,n_components,Y)
9 # -----
```

	cluster_0	cluster_1	cluster_2
setosa	50	0	0
versicolor	0	50	0
virginica	0	48	2

We observe a one-to-one correspondence between cluster 0 and the type iris setosa. The two other iris types (versicolor and virginica) are grouped into a single cluster. A last cluster contain two specific instances of the type virginica. These two instances appear to be somewhat distinct from other instances of the same type.

Exercise 4 (10 + 10 + 10 P)

In this exercise, we consider clustering of fashion items from the FashionMNIST dataset. The FashionMNIST dataset consists of 60000 fashion items, each of which coming as a 28×28 grayscale image. For the purpose of limiting computations, we consider a

subset of 1000 instances from this dataset, and thus, extract a dataset of size 1000 x 784 (the 784 dimensions correspond to representing images as a flat vector)

```
In [6]: 1 X = torchvision.datasets.FashionMNIST('.',download=True).data.numpy()
2
3 numpy.random.seed(0)
4 R = numpy.random.permutation(len(X))[:1000]
5 X = X[R].reshape(-1,784)
6
7 X = X / 255.0
```

The first 50 images of our dataset are visualized below:

```
In [7]: 1 plt.figure(figsize=(10,5))
2 plt.axis('off')
3 plt.imshow(X[:50].reshape(5,10,28,28).transpose(0,2,1,3).reshape(5*28,10*28),cmap='gray',vmin=
4 plt.show()
```

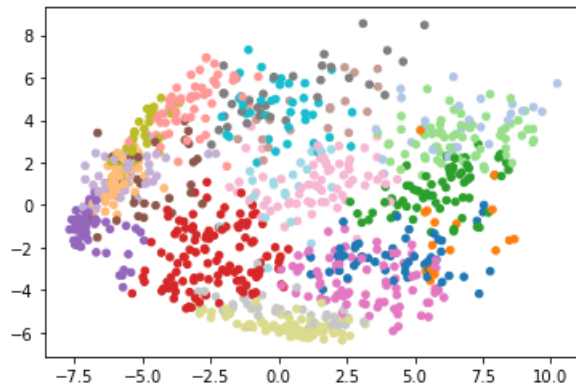


(a) Implement a rudimentary version of K-Means, where the cluster centers are initialized to the first few examples of the dataset, and where the number of iterations is fixed to 100.

```
In [8]: 1 import sklearn.cluster
2
3 class KMeans:
4
5     def __init__(self,n_clusters):
6         self.n_clusters = n_clusters
7
8     def fit(self,X):
9
10        # -----
11        # TODO: Replace by your code
12        # -----
13        import solution
14        solution.kmeans_fit(self,X)
15        # -----
16
17    def predict(self,X):
18
19        D = scipy.spatial.distance.cdist(X,self.cluster_centers_)
20        return numpy.argmin(D,axis=1)
21
22 model = KMeans(20)
23 model.fit(X)
```

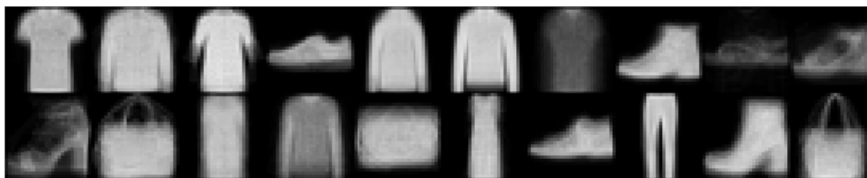
The outcome of your clustering procedure can be visualized in the following PCA plot color-coded by cluster membership.

```
In [9]: 1 numpy.random.seed(0)
2 Y = model.predict(X)
3 Z = sklearn.decomposition.PCA(n_components=2).fit_transform(X)
4 plt.scatter(*Z.T,c=Y,s=20,cmap='tab20')
5 plt.show()
```



(b) We now would like to get insights into what are prototypical fashion items in our dataset. For this, we consider the cluster centers learned by k-means. *Extract* these cluster centers and *visualize* them as images.

```
In [10]: 1 # -----
2 # TODO: Replace by your code
3 # -----
4 import solution
5 solution.view_cluster_centers(model)
6 # -----
```



(c) Lastly, we would like to know how well these cluster centers describe our data. Using the formulas presented in the lecture, compute the total variance, the within-cluster variance, the between-cluster variance, and the percentage of explained variance.

```
In [11]: 1 # -----
2 # TODO: Replace by your code
3 # -----
4 import solution
5 solution.analyze_variance(X,model)
6 # -----
```

Total variance	68.120
Between-cluster variance	41.579
Within-cluster variance	26.541
Explained variance (%)	61.04%

We observe that a bit more than half of the variance in the data is captured by the clustering, which is quite good, considering that the data is high-dimensional. Note that a higher explained variance could be obtained by incorporating more clusters, however that would come at the cost of making the data description potentially too complex for the end user.