

K-Means using scikit-learn

Use the scikit k-Means implementation to build the cluster of the data frame.

Preparations

Create the same data frame as above so that it is fresh.

```
In [4]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 import copy
6 import sklearn as sk
7 from sklearn.cluster import KMeans
8 # to check the time of execution, import function time
9 import time
10 # check versions of libraries
11 print('pandas version is: {}'.format(pd.__version__))
12 print('numpy version is: {}'.format(np.__version__))
13 print('sklearn version is: {}'.format(sk.__version__))
```

```
pandas version is: 1.0.1
numpy version is: 1.18.1
sklearn version is: 0.22.1
```

```
In [5]: 1 # Dataset
2 df = pd.DataFrame({
3     'x': [1, 2, 4, 5, 6, 8, 3, 7],
4     'y': [1, 1, 3, 4, 7, 8, 3, 7]
5 })
6
7 # Check that the definition of dataset is OK
8 print (""" data frame """)
9 print ("First column = No.")
10 print (df)
11
```

```
*** data frame ***
```

```
First column = No.
```

```
   x  y
0  1  1
1  2  1
2  4  3
3  5  4
4  6  7
5  8  8
6  3  3
7  7  7
```

K-Means training

Invoke the imported k-Means constructor with the number of clusters (here 3). Then train the model with the dataset.

```
In [6]: 1 # invoke constructor
2 kmeans = KMeans(n_clusters=3)
3
4 # Fitting K-Means model
5 print(kmeans.fit(df))
```

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=None, tol=0.0001, verbose=0)
```

K-Means prediction

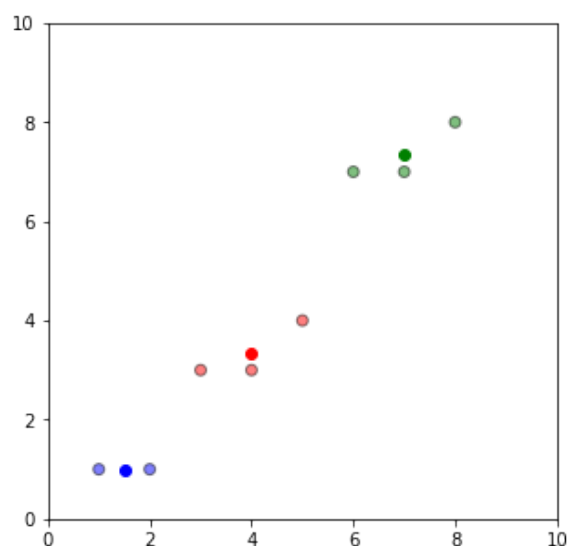
Use the model to calculate a prediction for the same data frame. Each datapoint will be labeled for the chosen cluster.

```
In [9]: 1 # create label for each datapoint in dataframe
2 labels = kmeans.predict(df)
3
4 # save centroids of each cluster
5 centroids = kmeans.cluster_centers_
```

Display the result

Display the positions of the centroids and the data frame. The color depends of the assigned label for each datapoint.

```
In [10]: 1 # Display result
2 fig = plt.figure(figsize=(5, 5))
3
4 # set color for each datapoint
5 colormap = {1: 'b', 2: 'g', 3: 'r'}
6 colors = list( map(lambda x: colormap[x+1], labels))
7
8 # draw each datapoint
9 plt.scatter(df['x'], df['y'], color=colors, alpha=0.5, edgecolor='k')
10
11 # draw each centroid
12 for idx, centroid in enumerate(centroids):
13     plt.scatter(*centroid, color=colormap[idx+1])
14 plt.xlim(0, 10)
15 plt.ylim(0, 10)
16 plt.show()
```



In [15]:

```
1
2 # print current date and time
3 print("date & time:",time.strftime("%d.%m.%Y %H:%M:%S"))
4 print ("*** End of Homework-H3.4_k-Means_Clustering ***")
5
```

date & time: 24.05.2023 22:31:59

*** End of Homework-H3.4_k-Means_Clustering ***