

Market segmentation example

In this notebook we explore a bit more sophisticated example of clustering

Import the relevant libraries

In []:



```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 # Set the styles to Seaborn
6 sns.set()
7 # Import the KMeans module so we can perform k-means clustering with sklearn
8 from sklearn.cluster import KMeans
```

Load the data

In []:



```
1 # Load the data
2 data = pd.read_csv ('3.12. Example.csv')
```

In []:



```
1 # Check what's inside
2 data
```

Plot the data

Create a preliminary plot to see if you can spot something

In []:



```
1 # We are creating a scatter plot of the two variables
2 plt.scatter(data['Satisfaction'],data['Loyalty'])
3 # Name your axes
4 plt.xlabel('Satisfaction')
5 plt.ylabel('Loyalty')
```

Select the features

In []:



```
1 # Select both features by creating a copy of the data variable
2 x = data.copy()
```

Clustering

In []:



```
1 # Create an object (which we would call kmeans)
2 # The number in the brackets is K, or the number of clusters we are aiming for
3 kmeans = KMeans(2)
4 # Fit the data
5 kmeans.fit(x)
```

Clustering results

In []:



```
1 # Create a copy of the input data
2 clusters = x.copy()
3 # Take note of the predicted clusters
4 clusters['cluster_pred']=kmeans.fit_predict(x)
```

In []:



```
1 # Plot the data using the Longitude and the Latitude
2 # c (color) is an argument which could be coded with a variable
3 # The variable in this case has values 0,1, indicating to plt.scatter, that there are 2 clusters
4 # All points in cluster 0 will be the same colour, all points in cluster 1 - another one
5 # cmap is the color map. Rainbow is a nice one, but you can check others here: https://matplotlib.org/1.2.2/mpl_toolkits/axes_helpers.html#color-maps
6 plt.scatter(clusters['Satisfaction'],clusters['Loyalty'],c=clusters['cluster_pred'],cmap=plt.cm.rainbow)
7 plt.xlabel('Satisfaction')
8 plt.ylabel('Loyalty')
```

Standardize the variables

Let's standardize and check the new result

In []:



```
1 # Import a library which can do that easily
2 from sklearn import preprocessing
3 # Scale the inputs
4 # preprocessing.scale scales each variable (column in x) with respect to itself
5 # The new result is an array
6 x_scaled = preprocessing.scale(x)
7 x_scaled
```

Take advantage of the Elbow method

In []:



```
1 # Createa an empty list
2 wcss =[]
3
4 # Create all possible cluster solutions with a loop
5 # We have chosen to get solutions from 1 to 9 clusters; you can ammend that if you wish
6 for i in range(1,10):
7     # Clsuter solution with i clusters
8     kmeans = KMeans(i)
9     # Fit the STANDARDIZED data
10    kmeans.fit(x_scaled)
11    # Append the WCSS for the iteration
12    wcss.append(kmeans.inertia_)
13
14 # Check the result
15 wcss
```

In []:



```
1 # Plot the number of clusters vs WCSS
2 plt.plot(range(1,10),wcss)
3 # Name your axes
4 plt.xlabel('Number of clusters')
5 plt.ylabel('WCSS')
```