

Assignment 2 - Clustering

Negin Baghbanzadeh

Learning Outcomes

In this assignment, you will do the following:

- Explore a dataset and carry out clustering using k-means algorithm
- Identify the optimum number of clusters for a given dataset

Problem Description

In this assignment, you will study the electricity demand from clients in Portugal, during 2013 and 2014. You have been provided with the data file, which you should download when you download this assignment file.

The data¹ available contains 370 time series, corresponding to the electric demand² for 370 clients, between 2011 and 2014.

In this guided exercise, you will use clustering techniques to understand the typical usage behaviour during 2013-2014.

Both these datasets are publicly available, and can be used to carry out experiments. Their source is below:

1. Data: <https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014#> (<https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014#>)
2. Electric Demand: <http://www.think-energy.net/KWvsKWH.htm> (<http://www.think-energy.net/KWvsKWH.htm>)

We will start by exploring the data set and continue on to the assignment. Consider this as a working notebook, you will add your work to the same notebook.

In this assignment we will use the sklearn package for k-means. Please refer here for the documentation <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html> (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>) (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>) (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>)).

The sklearn package for k-means is one of the many clustering algorithms found in the module "sklearn.cluster". These come with a variety of functions that you can call by importing the package.

For example

```
from sklearn.cluster import AgglomerativeClustering
from sklearn.cluster import KMeans
```

Data Preparation

Start by downloading the data to a local directory and modify the "pathToFile" and "fileName" variables, if needed. The data file has been provided with this assignment. It is also available at the links provided above.

```
In [1]: 1 pathToFile = r"/Users/Negin/Downloads/assignment2-data/"
        2
        3 fileName = 'LD2011_2014.txt'
```

```
In [2]: 1 import numpy as np
        2 from sklearn.cluster import KMeans
        3 import matplotlib.pyplot as plt
        4 import random
        5 from sklearn.metrics import silhouette_score
        6 from sklearn.cluster import AgglomerativeClustering
        7 random.seed(42)
```

```
In [3]: 1 # Replace "," by ".", otherwise the numbers will be in the form 2,
        2 import fileinput
        3
        4 with fileinput.FileInput(pathToFile+fileName, inplace=True, backup
        5     for line in file:
        6         print(line.replace(",", "."), end='')
```

```
In [4]: 1 # Create dataframe
        2 import pandas as pd
        3 data = pd.read_csv(pathToFile+fileName, sep=";", index_col=0)
```

Quick data inspection

In [5]: `1 data.head(2)`

Out [5]:

	MT_001	MT_002	MT_003	MT_004	MT_005	MT_006	MT_007	MT_008	MT_009	MT
2011-01-01 00:15:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
2011-01-01 00:30:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	

2 rows × 370 columns

In [6]: `1 data.tail(2)`

Out [6]:

	MT_001	MT_002	MT_003	MT_004	MT_005	MT_006	MT_007	MT_0
2014-12-31 23:45:00	1.269036	21.337127	1.737619	166.666667	85.365854	285.714286	10.17524	225.5892
2015-01-01 00:00:00	2.538071	19.914651	1.737619	178.861789	84.146341	279.761905	10.17524	249.1582

2 rows × 370 columns

In [7]: `1 data.shape`

Out [7]: (140256, 370)

As it can be seen, the dataframe contains a row for each interval of 15 minutes between Jan 1, 2011 to Dec 31 2014. There are 370 columns corresponding 370 clients. The dataframe is indexed by the timestamp.

Since the frequency is 15 minutes, each day provides $24 \times 4 = 96$ datapoints, which multiplied by 365 days and 4 years (plus 1 day in Feb 29, 2012) gives:
 $96 \times 365 \times 4 + 96 = 140256$, as observed in data.shape

In [8]: `1 data.info()`

```
<class 'pandas.core.frame.DataFrame'>
Index: 140256 entries, 2011-01-01 00:15:00 to 2015-01-01 00:00:00
Columns: 370 entries, MT_001 to MT_370
dtypes: float64(370)
memory usage: 397.0+ MB
```

In [9]: `1 data.describe()`

Out[9]:

	MT_001	MT_002	MT_003	MT_004	MT_005	MT_006
count	140256.000000	140256.000000	140256.000000	140256.000000	140256.000000	140256.000000
mean	3.970785	20.768480	2.918308	82.184490	37.240309	141.227000
std	5.983965	13.272415	11.014456	58.248392	26.461327	98.439000
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	2.844950	0.000000	36.585366	15.853659	71.428000
50%	1.269036	24.893314	1.737619	87.398374	39.024390	157.738000
75%	2.538071	29.871977	1.737619	115.853659	54.878049	205.357000
max	48.223350	115.220484	151.172893	321.138211	150.000000	535.714000

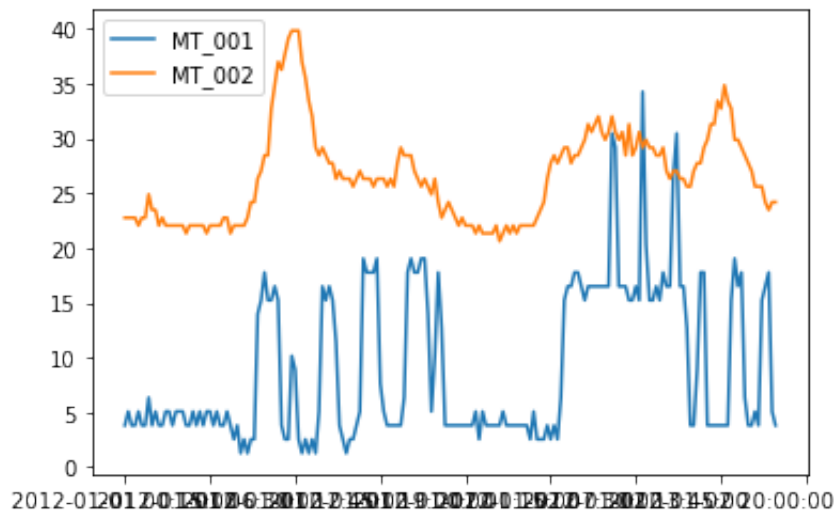
8 rows × 370 columns

In [10]: `1 data.isna().any().all()`

Out[10]: False

Plot the first 2 days of 2012 for the first 2 clients

In [11]: `1 data_example = data.loc['2012-01-01 00:15:00':'2012-01-03 00:00:00']
2 data_example.plot()
3 plt.show()`



We note that the main difference between the curves is the level (also seen on the means in `df.describe()`). We know we will have to somehow normalize the curves, in order for the clustering technique to capture the behaviour of the consumption throughout the day.

Data for the analysis

We focus on 2013 and 2014 because these are the years with low number of clients having zero demand

```
In [12]: 1 data2011 = data.loc['2011-01-01 00:15:00':'2012-01-01 00:00:00']
          2 data2012 = data.loc['2012-01-01 00:15:00':'2013-01-01 00:00:00']
          3 data2013 = data.loc['2013-01-01 00:15:00':'2014-01-01 00:00:00']
          4 data2014 = data.loc['2014-01-01 00:15:00':'2015-01-01 00:00:00']
```

```
In [13]: 1 # Check number of days
          2 print(data2011.shape[0]/96)
          3 print(data2012.shape[0]/96)
          4 print(data2013.shape[0]/96)
          5 print(data2014.shape[0]/96)
```

```
365.0
366.0
365.0
365.0
```

```
In [14]: 1 # See number of clients with 0 demand per year
          2 print(sum(data2011.mean()==0))
          3 print(sum(data2012.mean()==0))
          4 print(sum(data2013.mean()==0))
          5 print(sum(data2014.mean()==0))
```

```
210
37
21
1
```

```
In [15]: 1 clients = data2011.columns
2 clients_no_demand = clients[data2013.mean()==0] # clients with 0 d
3 data_13_14 = data2013.append(data2014) # appending 2013 and 2014
4 data_13_14 = data_13_14.drop(clients_no_demand, axis=1) # drop cli
5 print(data_13_14.shape)
6 print(sum(data_13_14.mean()==0)) # check that there are no clients
```

```
(70080, 349)
0
```

Getting average curves per client

```
In [16]: 1 data = data_13_14.copy() # weekdays weekends, data2011, data2012,
```

```
In [17]: 1 data['hour'] = data.index.map(lambda x: x[11:])
```

```
In [18]: 1 data.head(3)
```

Out[18]:

	MT_001	MT_002	MT_003	MT_004	MT_005	MT_006	MT_007	MT_0
2013-01-01 00:15:00	2.538071	22.759602	2.606429	138.211382	63.414634	255.952381	4.522329	239.0572
2013-01-01 00:30:00	1.269036	22.759602	2.606429	138.211382	63.414634	264.880952	5.652911	228.9562
2013-01-01 00:45:00	2.538071	22.759602	2.606429	134.146341	60.975610	250.000000	5.652911	239.0572

3 rows × 350 columns

Getting average curves per client

```
In [19]: 1 datagrouped = data.groupby("hour")
2 average_curves = datagrouped.agg("mean")
3 average_curves.shape
```

Out[19]: (96, 349)

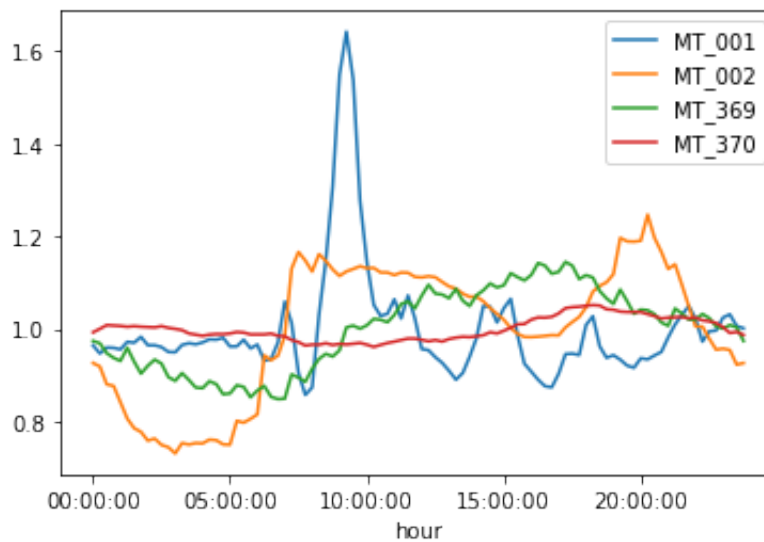
The dataframe `average_curves` contains the 349 typical weekday electric demands

We now divide each curve by its mean, so that all curves have mean 1

```
In [20]: 1 average_curves_norm = average_curves/(average_curves.mean())
```

Plot the first 2 and last 2 clients

```
In [21]: 1 average_curves_norm[['MT_001', 'MT_002', 'MT_369', 'MT_370']].plot()  
2 plt.show()
```



Clustering Analysis on the average normalized curves

```
In [22]: 1 X = average_curves_norm.copy() # We call this normalized curve  
2 X = np.array(X.T) # put it on the right format
```


Questions (15 marks total)

Q1: (7 marks)

a. Determine what a convenient number of clusters. Justify your choice. Make use of the sklearn's package for k-means for this. You may refer to the module to figure out how to come up with the optimal number of clusters.

b. Make a plot for each cluster, that includes: - The number of clients in the cluster (you can put this in the title of the plot) - All the curves in the cluster - The curve corresponding to the center of the cluster (make this curve thicker to distinguish it from the individual curves). The center is also sometimes referred to as "centroid".

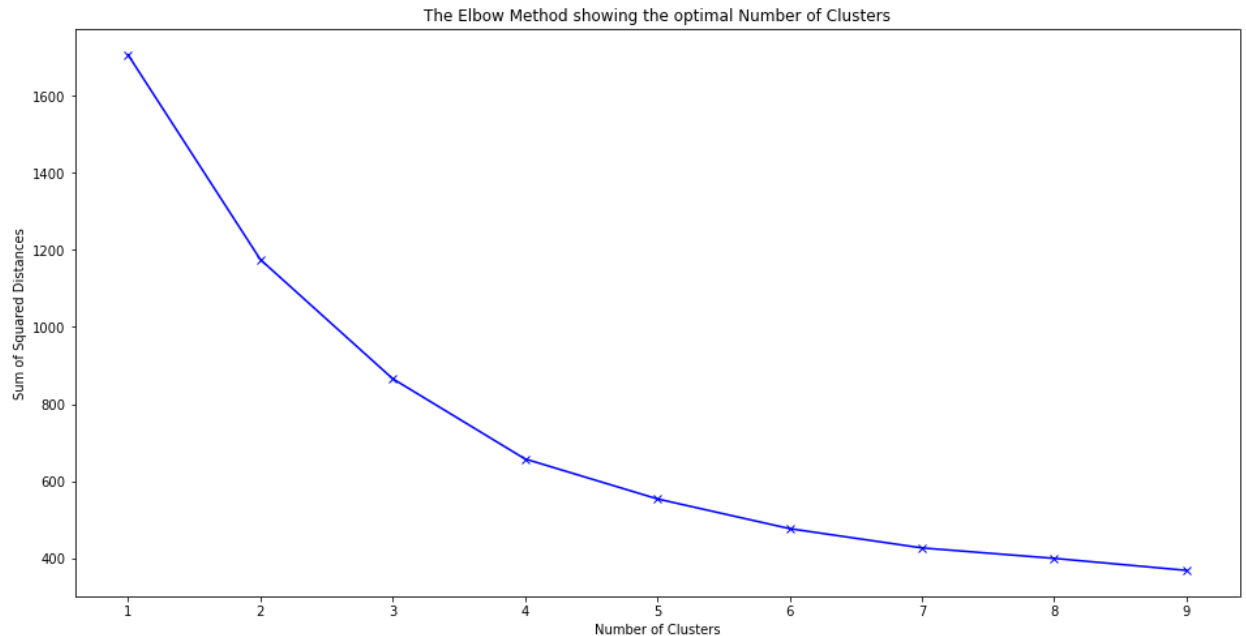
You have 2 separate plots for each cluster if you prefer (one for the individual curves, one for the centroid)

Q1 - Part a

In [23]:

```
1 sum_squared_dists = []
2 number_of_clusters = range(1,10)
3 for k in number_of_clusters:
4     kmeanModel = KMeans(n_clusters=k)
5     kmeanModel.fit(X)
6     sum_squared_dists.append(kmeanModel.inertia_)
```

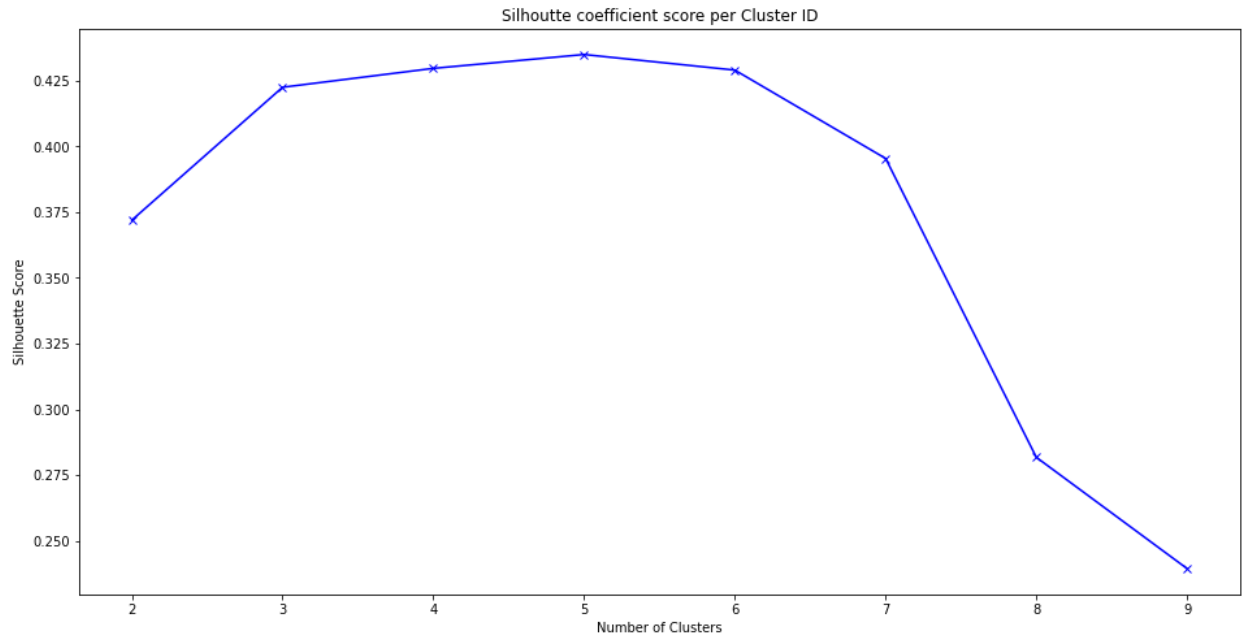
```
In [24]: 1 plt.figure(figsize=(16,8))
2 plt.plot(number_of_clusters, sum_squared_dists, 'bx-')
3 plt.xlabel('Number of Clusters')
4 plt.ylabel('Sum of Squared Distances')
5 plt.title('The Elbow Method showing the optimal Number of Clusters')
6 plt.show()
```



Based on the elbow graph above, having 6 clusters, doesn't make the much difference in making sum of squared distaances less than 7 clusters (making the model better). So a convenient number of clusters is **4** clusters.

```
In [25]: 1 from sklearn import metrics
2 silhouette_score = []
3 number_of_clusters = range(2,10)
4 for k in number_of_clusters:
5     kmeans_model_one_client = KMeans(n_clusters=k).fit(X)
6     labels = kmeans_model_one_client.labels_
7     silhouette_score.append(metrics.silhouette_score(X, labels, me
```

```
In [26]: 1 plt.figure(figsize=(16,8))
2 plt.plot(number_of_clusters, silhouette_score, 'bx-')
3 plt.xlabel('Number of Clusters')
4 plt.ylabel('Silhouette Score')
5 plt.title('Silhoutte coefficient score per Cluster ID')
6 plt.show()
```



The Silhoutte coefficient score confirm k=5 too.

Q1 - Part b

```
In [27]: 1 optimal_clusters = 5
2 optimal_model = KMeans(n_clusters=optimal_clusters).fit(X)
3 y_pred = optimal_model.labels_
```

```
In [28]: 1 from datetime import datetime
2
3 hours = average_curves_norm.index.map(lambda x: datetime.strptime(
```

```
In [29]: 1 # Finding the data in each cluster
2 clusters = []
3 for i in range(optimal_clusters):
4     clusters.append(X[y_pred == i])
```

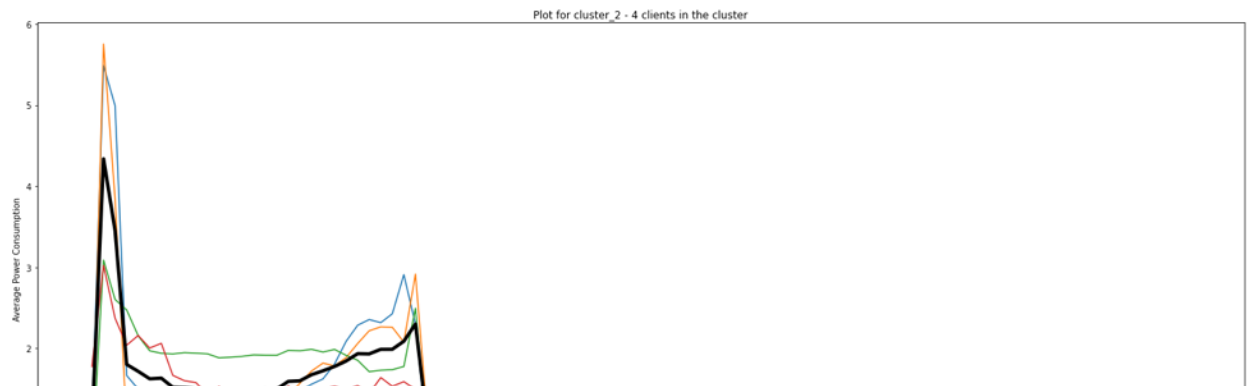
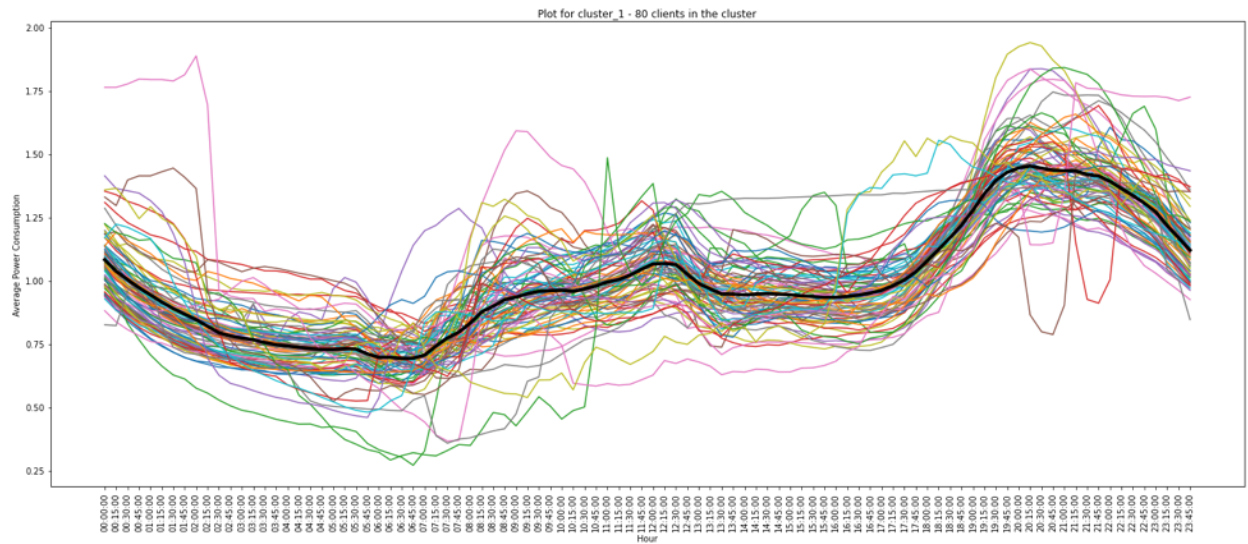
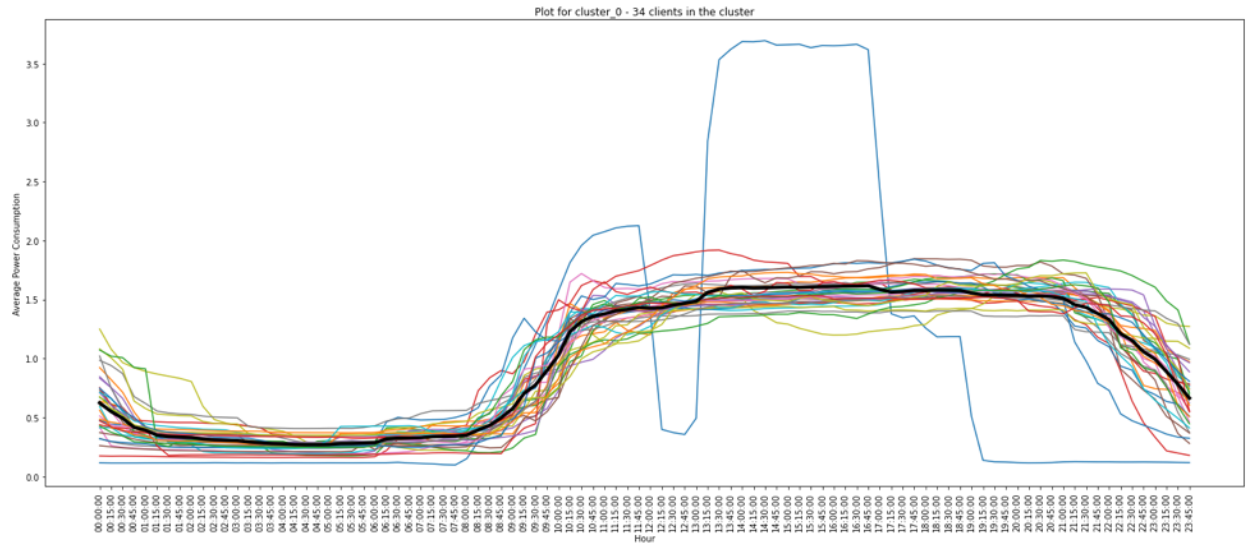
The centroids is shown by the black curve.

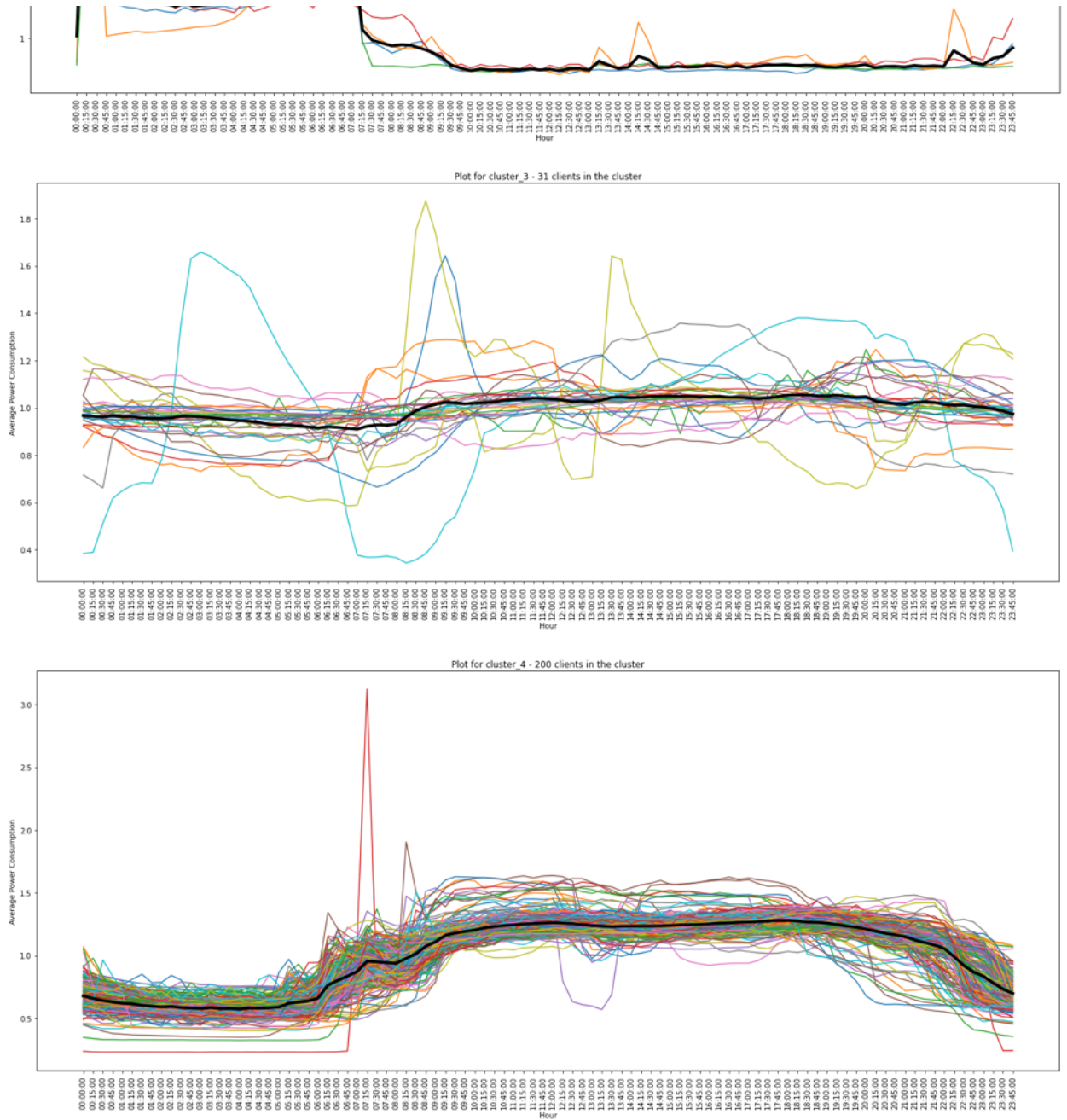
```
In [30]: 1 for i in range(optimal_clusters):
```

```

101 for i in range(optimal_clusters):
102     plt.figure(figsize=(25,10))
103     plt.xticks(rotation=90)
104     plt.plot(average_curves_norm.index, clusters[i].T)
105     plt.plot(average_curves_norm.index, optimal_model.cluster_centers[i].T)
106     plt.title("Plot for cluster_" + str(i) + " - " + str(len(clusters[i].T)))
107     plt.xlabel("Hour")
108     plt.ylabel("Average Power Consumption")

```





Q2: (8 marks)

In this exercise you work with the daily curves of 1 single client. First, create a list of arrays, each array containing a curve for a day. You may use X from the cells above. $X = \text{average_curves_norm.copy()}$ The list contains 730 arrays, one for each of the days of 2013 and 2014.

- Determine the optimal value of k (number of clusters). This time you may also perform silhouette analysis as stated in the module. Carrying out silhouette analysis is left as an exercise. What do you understand about the clusters?
- Based on your results from your analyses of both methods, what do understand? Interpret it perhaps with different perspectives of timelines like weeks or months.

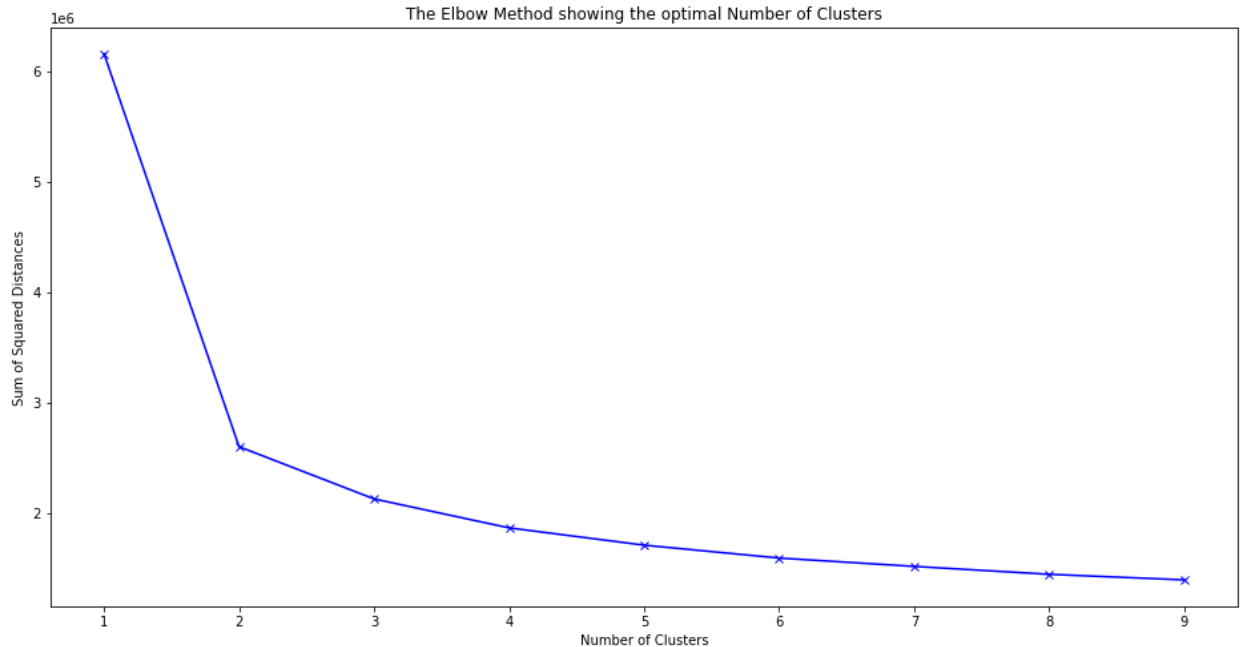
Q2 - Part a

```
In [31]: 1 client = 'MT_022'
          2 oneClient = data_13_14[client]
          3 X_one_client = [] # a list of arrays, each array being a normalized
          4 for J in range(2*365):
          5     X_one_client.extend([np.array(oneClient[J*96:(J+1)*96])]/np.
```

First I'm going to use Elbow Graph.

```
In [32]: 1 sum_squared_dists = []
          2 number_of_clusters = range(1,10)
          3 for k in number_of_clusters:
          4     kmeanModel = KMeans(n_clusters=k)
          5     kmeanModel.fit(X_one_client)
          6     sum_squared_dists.append(kmeanModel.inertia_)
```

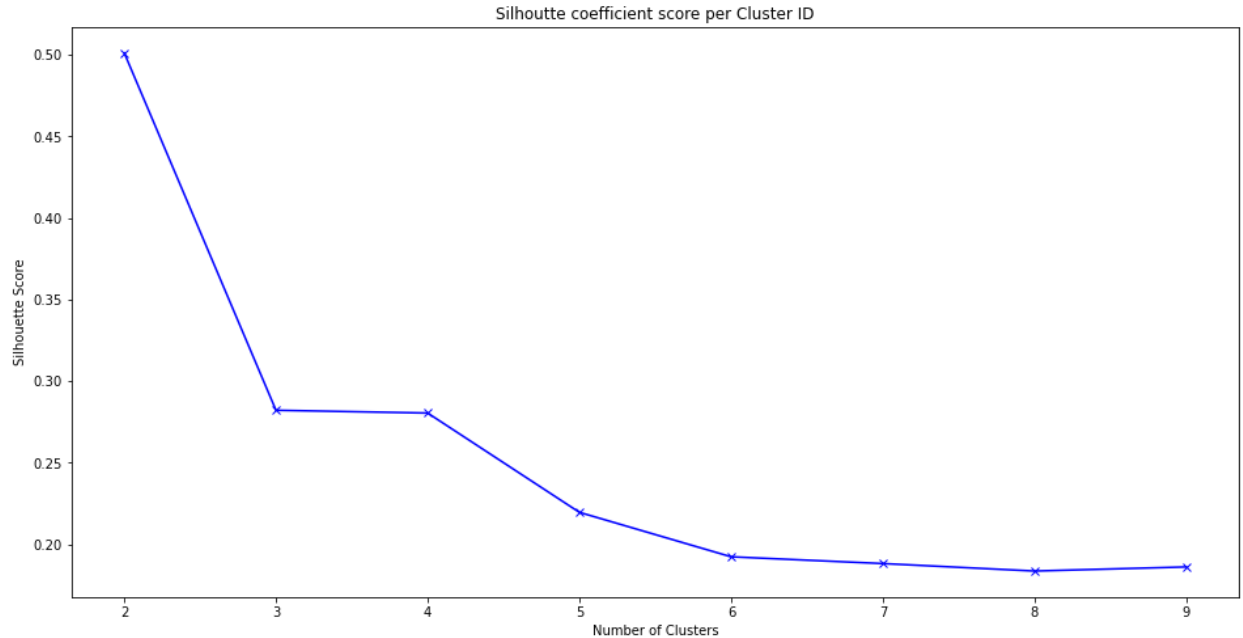
```
In [33]: 1 plt.figure(figsize=(16,8))
2 plt.plot(number_of_clusters, sum_squared_dists, 'bx-')
3 plt.xlabel('Number of Clusters')
4 plt.ylabel('Sum of Squared Distances')
5 plt.title('The Elbow Method showing the optimal Number of Clusters')
6 plt.show()
```



Based on the graph, I think k=3 is a good number for the clusters. (but k=2 also seems like good number of clusters)

```
In [34]: 1 silhouette_score = []
2 number_of_clusters = range(2,10)
3 for k in number_of_clusters:
4     kmeans_model_one_client = KMeans(n_clusters=k).fit(X_one_client)
5     labels = kmeans_model_one_client.labels_
6     silhouette_score.append(metrics.silhouette_score(X_one_client,
```

```
In [35]: 1 plt.figure(figsize=(16,8))
2 plt.plot(number_of_clusters, silhouette_score, 'bx-')
3 plt.xlabel('Number of Clusters')
4 plt.ylabel('Silhouette Score')
5 plt.title('Silhoutte coefficient score per Cluster ID')
6 plt.show()
```



$s(x)$ is a quantity in $[-1, 1]$. A larger Silhoutte coefficient score relates to a model with better defined clusters. Based on the graph, the best number of clusters is 2 clusters.

In the Elbow graph, I had doubt between 2 and 3(Although the change from 1 to 2 was way more than the change between 2 and 3), but based on the Silhoutte coefficient score, $k=2,3,4$ are better choices than the rest but $k=2$ is way better than $k=3$ and $k=4$. So we Can say that these two methodds confirm each other.

Q2 - Part b

The following cell generates a list of all the days in the two years 2013-2014, which is helpful to answer part (b). The cells below are only to give you a headstart. You may or may not use these and come up with your own interpretation.


```
In [36]: 1 from datetime import date, timedelta
2
3 d1 = date(2013, 1, 1) # start date
4 d2 = date(2014, 12, 31) # end date
5 delta = d2 - d1 # timedelta
6 daysyear = []
7 D = {0:'mon', 1:'tue', 2:'wed', 3:'thu', 4:'fri', 5:'sat', 6:'sun'}
8 for i in range(delta.days + 1):
9     daysyear.extend([D[(d1 + timedelta(days=i)).weekday()]]+"-"+str
```

Since we have 2 clusters, one easy guess is that the clusters represent weekdays and weekends. So we're going to plot days of the year vs cluster ID to see if our alternative hypothesis is right or not.

```
In [37]: 1 X_one_client_df = pd.DataFrame(X_one_client)
2 X_one_client_df['daysyear'] = daysyear #[i[0:3] for i in daysyear]
3 optimal_k=2
4 optimal_model_one_client = KMeans(n_clusters=optimal_k).fit(X_one_
5 X_one_client_df['cluster'] = optimal_model_one_client.labels_
```

```
In [38]: 1 plt.figure(figsize=(500,10))
2 plt.xticks(rotation=90)
3 plt.plot(X_one_client_df['daysyear'], X_one_client_df['cluster'])
4 plt.title("Days of Year VS. Cluster ID")
5 plt.xlabel("Days in Year")
6 plt.ylabel("Cluster ID")
```

```
Out[38]: Text(0, 0.5, 'Cluster ID')
```



If you look at the above graph(double click on it) you can see that most of the times, on weekends, the cluster ID is 0 and on weekdays it is 1. So we can say our alternative hypothesis is right and the clusters represent week days and weekends.