

Building our own classification model

In this notebook, we will create a classification model to help Perpetual Investors get insightful environmental and social information about the pilot companies. The goal is to help them make three different decisions about whether they should invest in these companies. The classification model will assign the companies in data three different categories regarding their Environmental Actions:

- 1) Good Forecasting 2) Improving 3) Bad forecasting

Further on the notebook, we will get into more detail about the required 'rules' for each classification group.

As mentioned, we are doing this classification model for the 52 US pilot stocks. Moreover, because we know how important the industry factor is, we decided to classify the companies within each industry to make them comparable. Therefore, we will have two classification model. One to classify the companies in the Utility industry and another for the Energy industry.

Steps followed to build our own classification model:

- 1) Data Cleaning 2) Hierarchical Clustering 3) K-Means Clustering 4) Cluster Profiling 5) Assigned the labels based in the final clusters and rules

We will have a Hierarchical Clustering and K-means clustering for each industry. We used Unsupervised Machine Learning to find patterns and similarities within the companies.

Let's start building the model.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

# For distance and h-clustering
from scipy.cluster.hierarchy import linkage, dendrogram, fcluster
from scipy.spatial.distance import pdist, squareform

# Means clustering
from sklearn.cluster import KMeans, DBSCAN
from sklearn.neighbors import NearestNeighbors
from sklearn.metrics import pairwise_distances
from sklearn.decomposition import PCA

# sklearn models have some functionality too, but mostly a wrapper to scipy
from sklearn.preprocessing import StandardScaler
import warnings
warnings.filterwarnings('ignore')

df = pd.read_csv('df_merged.csv')
df['cl'] = pd.read_csv('Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/Dataset/stocks')
stocks = pd.read_csv('Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/Dataset/stocks')

In [2]: sectors = pd.read_csv('Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/Dataset/stocks')['Missing_GHG']
stocks['Missing_GHG'] = np.where(stocks['GHG_Scope 1'].isna(), 1, 0)
df['df_merged(Stocks, sectors, how="inner", on="Ticker")']
df['df_merged(Stocks, sectors, how="inner", on="Ticker")']
df['df_merged(Stocks, sectors, how="inner", on="Ticker")']
stocks = df.copy()
stocks['Utility'] = np.where(stocks.Sector == 'Utilities', 1, 0)
stocks['Energy'] = 1 - stocks['Utility']

Out [3]: <class 'pandas.core.frame.DataFrame'>
Int64Index: 780 entries, 0 to 779
Data columns (total 28 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Year                 780 non-null    int64
1   Ticker               780 non-null    object
2   Environmental Disclosure Score  636 non-null    float64
3   GHG Scope 1          398 non-null    float64
4   Total Energy Consumption  243 non-null    float64
5   Change_in_EDS        589 non-null    float64
6   Change_in_GHG        237 non-null    float64
7   Change_in_NET        780 non-null    object
8   Company              780 non-null    object
9   Net_Income           780 non-null    float64
10  Total_Assets          754 non-null    float64
11  Total_Sales           754 non-null    float64
12  Change_in_Sales      754 non-null    float64
13  Change_in_Assets     754 non-null    float64
14  Change_in_NI         754 non-null    float64
15  ROA                  754 non-null    float64
16  Profit_Margin        754 non-null    float64
17  Annual_Stock_Return  366 non-null    float64
18  Ratio of GHG Emissions to Total Assets  398 non-null    float64
19  Ratio of GHG Emissions to Total Sales  398 non-null    float64
20  Ratio of Total Energy Consumption to Total Assets  243 non-null    float64
21  Ratio of Total Energy Consumption to Total Sales  243 non-null    float64
22  Profitable           780 non-null    bool
23  Logarithm_Total_Assets  754 non-null    float64
24  Logarithm_Total_Sales  754 non-null    float64
25  Utility              780 non-null    bool
26  Sector               780 non-null    object
27  Utility              780 non-null    int64
memory usage: 171.4+ KB
```

Data Cleaning - version 1

As presented throughout the project, industry plays a crucial factor when trying to understand, interpret and predict the company environmental actions. In this version, we decided to impute the missing values using the industry average of that year.

At the end, we decided to not use this version. However, we want to present that it was part of our analysis.

```
In [3]: stocks_df = stocks_df.copy()
stocks_df = stocks_df.loc[:, ~stocks_df.columns.isin(['Logarithm_Total_Assets', 'Logarithm_Total_Sales', 'Sector', 'Logarithm_Total_Assets', 'Logarithm_Total_Sales', 'Sector', 'Logarithm_Total_Assets', 'Logarithm_Total_Sales', 'Sector', 'Logarithm_Total_Assets', 'Logarithm_Total_Sales', 'Sector'])]
stocks_df['df_merged(Stocks, sectors, how="inner", on="Ticker")'] = np.nan # We can see infinity value

In [4]: stocks_df.describe().T

Out [4]:
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|---|-------|---------------|--------------|---------------|--------------|---------------|---------------|---------------|
| Year | 780.0 | 2012.000000 | 4.22366 | 2005.000000 | 2008.000000 | 2012.000000 | 2016.000000 | 2.01900e+03 |
| Environmental Disclosure Score | 636.0 | 31.469954 | 18.263954 | 1.400000 | 15.900000 | 31.750000 | 45.500000 | 8.430000e+01 |
| GHG Scope 1 | 398.0 | 28308.54774 | 3206.231283 | 52.500000 | 5225.000000 | 45000.000000 | 36650.000000 | 1.455000e+05 |
| Total Energy Consumption | 243.0 | 138491.043210 | 30823.410807 | 101.500000 | 10630.400000 | 191643.500000 | 191643.500000 | 3.773300e+06 |
| Change_in_EDS | 589.0 | 25.285433 | 105.813056 | -65.363128 | 0.000000 | 2.725724 | 19.090099 | 1.887500e+03 |
| Change_in_GHG | 377.0 | -0.660034 | 26.111899 | -98.965689 | -8.586678 | -1.145310 | 3.196347 | 2.408630e+02 |
| Change_in_NET | 237.0 | 122.143943 | 1228.707439 | -98.965689 | -8.586678 | 0.000000 | 3.196347 | 1.575020e+05 |
| Total Assets | 754.0 | 44006.135068 | 5246.876654 | 172.517000 | 17640.800500 | 30182.100000 | 47309.500000 | 3.625970e+04 |
| Total Energy Consumption | 243.0 | 1862.769263 | 5374.043089 | -23119.000000 | 306.017050 | 827.396500 | 1830.750000 | 4.522000e+04 |
| Net Income | 754.0 | 27927.763937 | 56598.313351 | 27.253000 | 6639.350000 | 11485.987500 | 16794.000000 | 4.335260e+05 |
| Total Sales | 754.0 | 27497.763937 | 56598.313351 | 27.253000 | 6639.350000 | 11485.987500 | 16794.000000 | 4.335260e+05 |
| Change_in_Sales | 754.0 | 4.887176 | 319.94785 | -78.151774 | -5.295453 | 1.881040 | 11.135700 | 6.470430e+02 |
| Change_in_Assets | 754.0 | 7.192162 | 2.1082076 | -66.324707 | 0.000000 | 4.803945 | 9.095253 | 1.786625e+02 |
| Change_in_NI | 754.0 | -10.272885 | 926.832093 | -8982.500000 | -33.633989 | 0.021741 | 24.550031 | 1.956120e+04 |
| ROA | 754.0 | 3.092603 | 7.516937 | -122.699289 | 1.917490 | 2.855311 | 4.719352 | 3.541617e+01 |
| Profit Margin | 754.0 | 6.290280 | 20.66235 | -36.196459 | 3.837708 | 8.414402 | 12.040351 | 4.675778e+02 |
| Annual Stock Return | 362.0 | 1.945314 | 24.569376 | -62.309948 | -2.450866 | 11.257979 | 21.377972 | 1.3772194e+01 |
| Ratio of GHG Emissions to Total Assets | 398.0 | 590.490234 | 625.325587 | 3.041010 | 158.325338 | 363.467683 | 829.690452 | 4.677381e+03 |
| Ratio of GHG Emissions to Total Sales | 398.0 | 1838.463042 | 221.2110641 | 15.105900 | 282.722003 | 654.607648 | 3313.797971 | 8.998950e+03 |
| Ratio of Total Energy Consumption to Total Assets | 243.0 | 1928.396940 | 4716.675535 | 4.694293 | 312.373762 | 1027.323793 | 2140.752073 | 6.688640e+04 |
| Ratio of Total Energy Consumption to Total Sales | 243.0 | 5655.626268 | 16864.951584 | 15.861853 | 742.879876 | 1518.557341 | 6102.365655 | 2.457061e+05 |
| Utility | 780.0 | 0.538462 | 0.498838 | 0.000000 | 0.000000 | 1.000000 | 1.000000 | 1.000000e+00 |

Let's impute null values with that year industry average for all the columns that contain null values.

```
In [5]: null_columns = stocks_df.columns[stocks_df.isnull().any()]
null_columns = null_columns.tolist()

for column in null_columns:
    stocks_df[column] = stocks_df.groupby(['Year', 'Utility']).transform('mean')[column]

In [6]: stocks_df.isna().sum()

Out [6]:
```

| Year | 0 |
|---|-----|
| Year | 0 |
| Ticker | 0 |
| Environmental Disclosure Score | 0 |
| GHG Scope 1 | 28 |
| Total Energy Consumption | 56 |
| Change_in_EDS | 80 |
| Change_in_GHG | 108 |
| Company | 0 |
| Total Assets | 0 |
| Total Sales | 0 |
| Change_in_Sales | 0 |
| Change_in_Assets | 0 |
| Change_in_NI | 0 |
| ROA | 0 |
| Profit Margin | 0 |
| Annual Stock Return | 41 |
| Ratio of GHG Emissions to Total Assets | 28 |
| Ratio of GHG Emissions to Total Sales | 28 |
| Ratio of Total Energy Consumption to Total Assets | 56 |
| Ratio of Total Energy Consumption to Total Sales | 56 |
| Profitable | 0 |
| Utility | 0 |
| dtype: int64 | 0 |

We have less null values now. However, we keep seeing a great amount of null values for the Annual Stock Return and Change in Total Energy Consumption. We are going to proceed to drop these columns and not consider them for the clustering analysis.

```
In [7]: stocks_df = stocks_df.loc[:, ~stocks_df.columns.isin(['Annual_Stock_Return', 'Change_in_NET'])]
stocks_df.head()

Out [7]:
```

| Year | Ticker | Environmental Disclosure Score | GHG Scope 1 | Total Energy Consumption | Change_in_EDS | Change_in_GHG | Company | Total Assets | Net Income | ... | Change |
|------|--------|--------------------------------|-------------|--------------------------|---------------|---------------|-------------|--------------|--------------|------------|--------|
| 0 | 2005 | AEE | 16.566667 | NaN | NaN | NaN | AMEREN CORP | 24071.876643 | 576.930357 | .. | |
| 1 | 2006 | AEE | 19.153846 | 23845.733333 | NaN | 33.294759 | AMEREN CORP | 25461.556893 | 743.487464 | .. | |
| 2 | 2007 | AEE | 18.677273 | 29858.580000 | 342779.0 | 49.676880 | 2.352400 | AMEREN CORP | 25310.003393 | 849.919000 | .. |
| 3 | 2008 | AEE | 19.676923 | 23344.828571 | 327939.0 | 14.110217 | -1.376797 | AMEREN CORP | 27335.062393 | 897.238429 | .. |
| 4 | 2009 | AEE | 26.074074 | 18867.840000 | 164566.8 | 17.994489 | -8.848381 | AMEREN CORP | 28166.217107 | 793.816214 | .. |

5 rows x 22 columns

```
In [8]: is_NaN = stocks_df.isnull()
row_NaN = stocks_df[is_NaN.any(axis=1)]
rows_with_NaN = stocks_df[rows_NaN]
null_columns = rows_with_NaN.columns[rows_with_NaN.isnull().any()]
null_columns = null_columns.tolist()
null_columns.append('Year')
rows_with_NaN[null_columns]
```

| GHG Scope 1 | Total Energy Consumption | Change_in_EDS | Change_in_GHG | Ratio of GHG Emissions to Total Assets | Ratio of GHG Emissions to Total Sales | Ratio of Total Energy Consumption to Total Assets | Ratio of Total Energy Consumption to Total Sales | Year |
|-------------|--------------------------|---------------|---------------|--|---------------------------------------|---|--|-----------------|
| 0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 2005 |
| 1 | 23845.733333 | NaN | 33.294759 | NaN | 3352.017617 | NaN | NaN | 2006 |
| 15 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 2005 |
| 16 | 23845.733333 | NaN | 33.294759 | NaN | 3352.017617 | NaN | NaN | 2006 |
| 30 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 2005 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 721 | 23845.733333 | NaN | 33.294759 | NaN | 3352.017617 | NaN | NaN | 2006 |
| 735 | 11950.000000 | 74970.0 | NaN | NaN | 432.979753 | 410.866467 | 2630.710927 | 1279.43887 2005 |
| 750 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 2006 |
| 751 | 23845.733333 | NaN | 33.294759 | NaN | 3352.017617 | NaN | NaN | 2006 |
| 765 | 11950.000000 | 74970.0 | NaN | NaN | 432.979753 | 410.866467 | 2630.710927 | 1279.43887 2005 |

80 rows x 9 columns

```
In [9]: rows_with_NaN[null_columns] = rows_with_NaN.isnull().sum(axis=1)
rows_with_NaN.head()

plt.figure(figsize=(10,5))
sns.barplot(x='Year', y='null_values', data=rows_with_NaN)
```

Let's get rid of data from 2005 and 2006 because it is missing data for multiple columns

```
In [10]: stocks_df = stocks_df.loc[stocks_df.Year >= 2007, :]
stocks_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 676 entries, 2 to 779
Data columns (total 22 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Year                 676 non-null    int64
1   Ticker               676 non-null    object
2   Environmental Disclosure Score  676 non-null    float64
3   GHG Scope 1          676 non-null    float64
4   Total Energy Consumption  676 non-null    float64
5   Change_in_EDS        676 non-null    float64
6   Change_in_GHG        676 non-null    float64
7   Company              676 non-null    object
8   Total_Assets          676 non-null    float64
9   Net_Income           676 non-null    float64
10  Total_Sales           676 non-null    float64
11  Change_in_Sales      676 non-null    float64
12  Change_in_Assets     676 non-null    float64
13  Change_in_NI         676 non-null    float64
14  ROA                  676 non-null    float64
15  Profit_Margin        676 non-null    float64
16  Ratio of GHG Emissions to Total Assets  676 non-null    float64
17  Ratio of GHG Emissions to Total Sales  676 non-null    float64
18  Ratio of Total Energy Consumption to Total Assets  676 non-null    float64
19  Ratio of Total Energy Consumption to Total Sales  676 non-null    float64
20  Profitable           676 non-null    bool
21  Utility              676 non-null    bool
dtypes: bool(2), float64(17), int64(2), object(2)
memory usage: 116.8+ KB
```

Hierarchical Clustering - version 1 and utility industry

Now that we have a cleaned dataset, we will start with the clustering. First, we did Hierarchical Clustering for the utility industry.

As a reminder, we are not using this version for the classification model.

```
In [11]: # 1) Standardized the data
df = stocks_df[stocks_df.Utility == 1]
stock_number = df.select_dtypes('number')
sc = StandardScaler()
stock_scaled = sc.fit_transform(stock_number)
stock_scaled = pd.DataFrame(stock_scaled, columns = stock_number.columns)
stock_scaled

Out [11]:
```

| Year | Environmental Disclosure Score | GHG | Total Energy Consumption | Change_in_EDS | Change_in_GHG | Total Assets | Net Income | Total Sales | Change_in_Sale | |
|------|--------------------------------|-----------|--------------------------|---------------|---------------|--------------|------------|-------------|----------------|-----------|
| 0 | -1.603567 | -1.265688 | -0.099899 | 0.715048 | 0.879926 | 0.668179 | -1.458084 | -0.330097 | -1.326754 | 0.159041 |
| 1 | -1.136306 | -1.514224 | -1.193930 | 0.629004 | -0.332286 | 0.104514 | -1.222410 | -0.136021 | -0.150339 | 1.239266 |
| 2 | -1.069045 | -0.795163 | -2.007669 | -0.118249 | 3.208382 | -1.024811 | -1.125682 | -0.560196 | -1.284446 | -2.124061 |
| 3 | -0.801784 | -0.724807 | -1.137553 | -0.395566 | -0.466368 | 0.551322 | -0.925459 | -0.185963 | -1.045936 | 0.495355 |
| 4 | -0.534522 | -0.552455 | 0.729316 | 0.091868 | -0.289423 | -0.024700 | -0.641432 | 0.017768 | -0.567129 | 0.234044 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 359 | 0.354522 | 0.528622 | 0.776012 | -0.483558 | -0.065203 | -1.027699 | 0.385752 | -0.817127 | 0.425876 | -1.603188 |
| 730 | 0.801784 | 0.877982 | 0.295243 | -0.648703 | -0.146088 | -0.936210 | 0.808715 | -1.260166 | 0.253713 | -0.544833 |
| 361 | 1.069045 | 1.241836 | 0.314184 | -0.734232 | -0.264208 | -0.261941 | 0.972987 | 0.291251 | 0.707892 | 0.079317 |
| 362 | 1.363606 | 1.275571 | 0.025056 | -0.657463 | -0.7113375 | -0.249911 | 1.285929 | 1.425912 | 1.261793 | 0.717635 |
| 751 | 1.603567 | 1.507336 | 0.410070 | -0.704496 | -0.480345 | -1.032557 | 1.830051 | 2.622920 | 3.131953 | -0.718620 |

364 rows x 19 columns

```
In [12]: #2) Clustering using euclidean and cosine for distance matrix
dcl = pdist(stock_scaled.values) #euclidean
dcl = pdist(stock_scaled.values, metric='cosine')

#See how with linkage method and cosine distance matrix work
METHODS = ['single', 'complete', 'average', 'ward']
plt.figure(figsize=(20,5))

for i, m in enumerate(METHODS):
    plt.subplot(1,4,i+1)
    plt.title(m)
    dendrogram(linkage(dcl, method=m),
                leaf_rotation=90)
```

I am going to use the cosine and average because the cluster start forming lower and there isn't too much height compared to the other methods. Clusters are more group together and compact

```
In [14]: #4) Create the labels
hcl = linkage(dcl, method='average')
plt.title('Dendrogram for Cosine and Average')
dendrogram(hcl,
            leaf_rotation=90)
plt.axhline(linestyle='--', y=83)

Out [14]:
```

```
In [15]: labels = fcluster(hcl, 3, criterion='maxclust')
np.unique(labels)

Out [15]: array([1, 2, 3], dtype=int32)

In [16]: #put the labels into the clean dataset
df['cluster'] = labels
df.head(3) #Review the dataset with the labels

Out [16]:
```

| Year | Ticker | Environmental Disclosure Score | GHG Scope 1 | Total Energy Consumption | Change_in_EDS | Change_in_GHG | Company | Total Assets | Net Income | ... | Change |
|------|--------|--------------------------------|-------------|--------------------------|---------------|---------------|-----------|--------------|--------------|------------|--------|
| 2 | 2007 | AEE | 18.677273 | 29858.580000 | 342779.0 | 49.676880 | 2.352400 | AMEREN CORP | 25310.003393 | 849.919000 | .. |
| 3 | 2008 | AEE | 19.676923 | 23344.828571 | 327939.0 | 14.110217 | -1.376797 | AMEREN CORP | 27335.062393 | 897.238429 | .. |
| 4 | 2009 | AEE | 26.074074 | 18867.840000 | 164566.8 | 17.994489 | -8.848381 | AMEREN CORP | 28166.217107 | 793.816214 | .. |

3 rows x 23 columns

```
In [17]: #How many stocks per cluster
df1.cluster.value_counts(dropna=False, sort=False)

Out [17]:
```

| 1 | 168 |
|-----------------------------|-----|
| 2 | 112 |
| 3 | 84 |
| Name: cluster, dtype: int64 | |

When using three clusters... we can see that the number of observations are more balanced.

```
In [18]: X = stock_scaled.values
sns.scatterplot(X[:,3],X[:,4],hue=df1.cluster, cmap='rainbow').set(title='Stock - Hierarchical Clustering')

Out [18]:
```

K-Means Clustering - version 1 and utility

```
In [19]: #Cluster Evaluation - Deciding how many clusters
X = stock_scaled.values
KRANGE = range(2,10)
sse = []

# Loop over and evaluate
for k in KRANGE:
    km = KMeans(k)
    lab = km.fit_predict(stock_scaled)
    sse.append(km.inertia_)

#Elbow Method
sns.lineplot(KRANGE,sse)

Out [19]:
```

```
In [20]: # Testing K
ksl = []

for k in KRANGE:
    km = KMeans(k)
    lab = km.fit_predict(stock_scaled)
    sse.append(metrics.silhouette_score(stock_scaled, lab))

sns.lineplot(KRANGE, sse)

Out [20]:
```

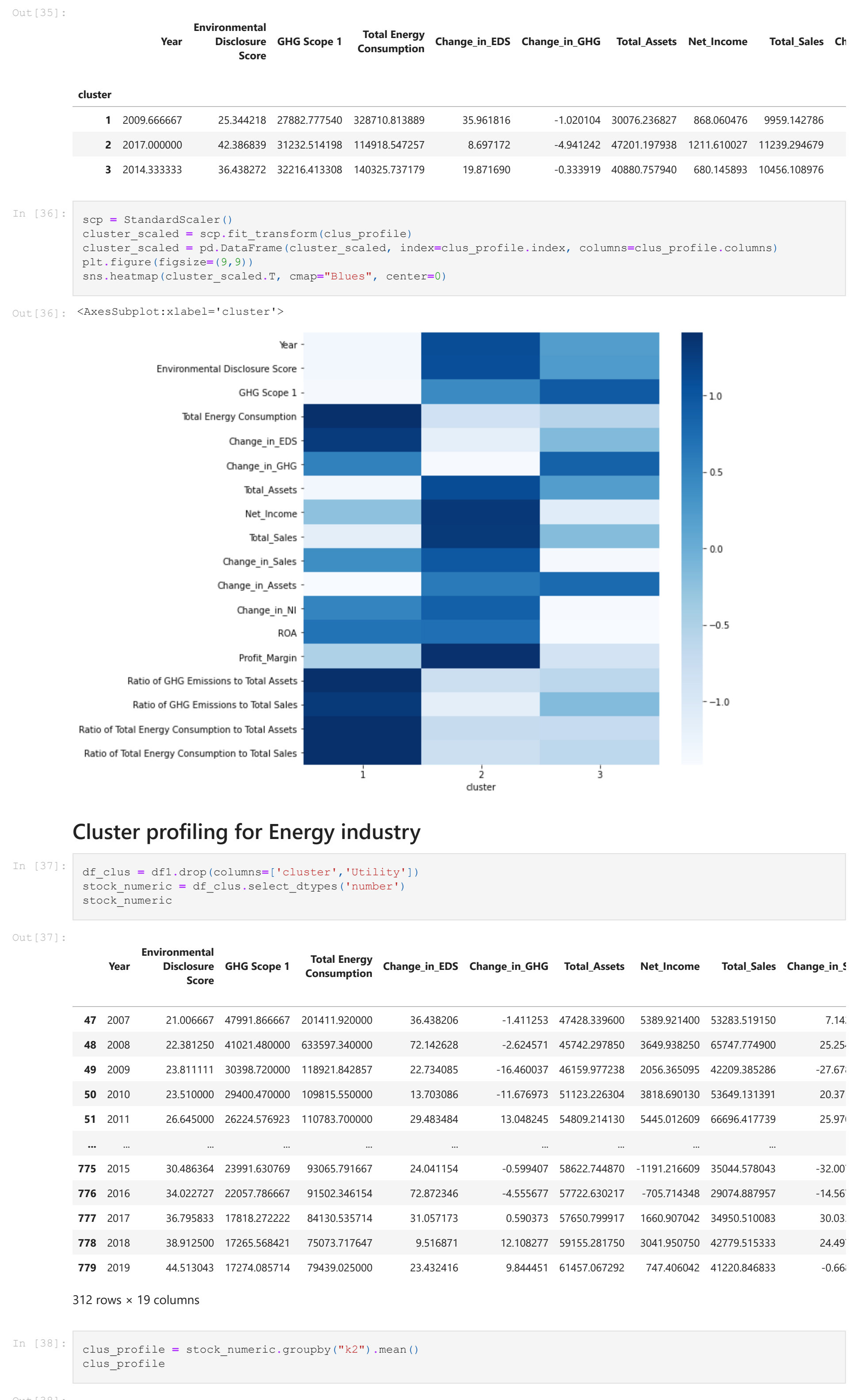
```
In [21]: #It shows 3 clusters are the best option... because it radically decreases when 4.

##K Means for 3
k2 = KMeans(3)
k2.fit(X)
k2_labels = k2.predict(X)

In [22]: df['k3'] = k3_labels
df.k3.value_counts(dropna=False, sort=False) #

Out [22]:
```

| 0 | 84 |
|-------|-----|
| 1 | 112 |
| 2 | 168 |
| Name: | |



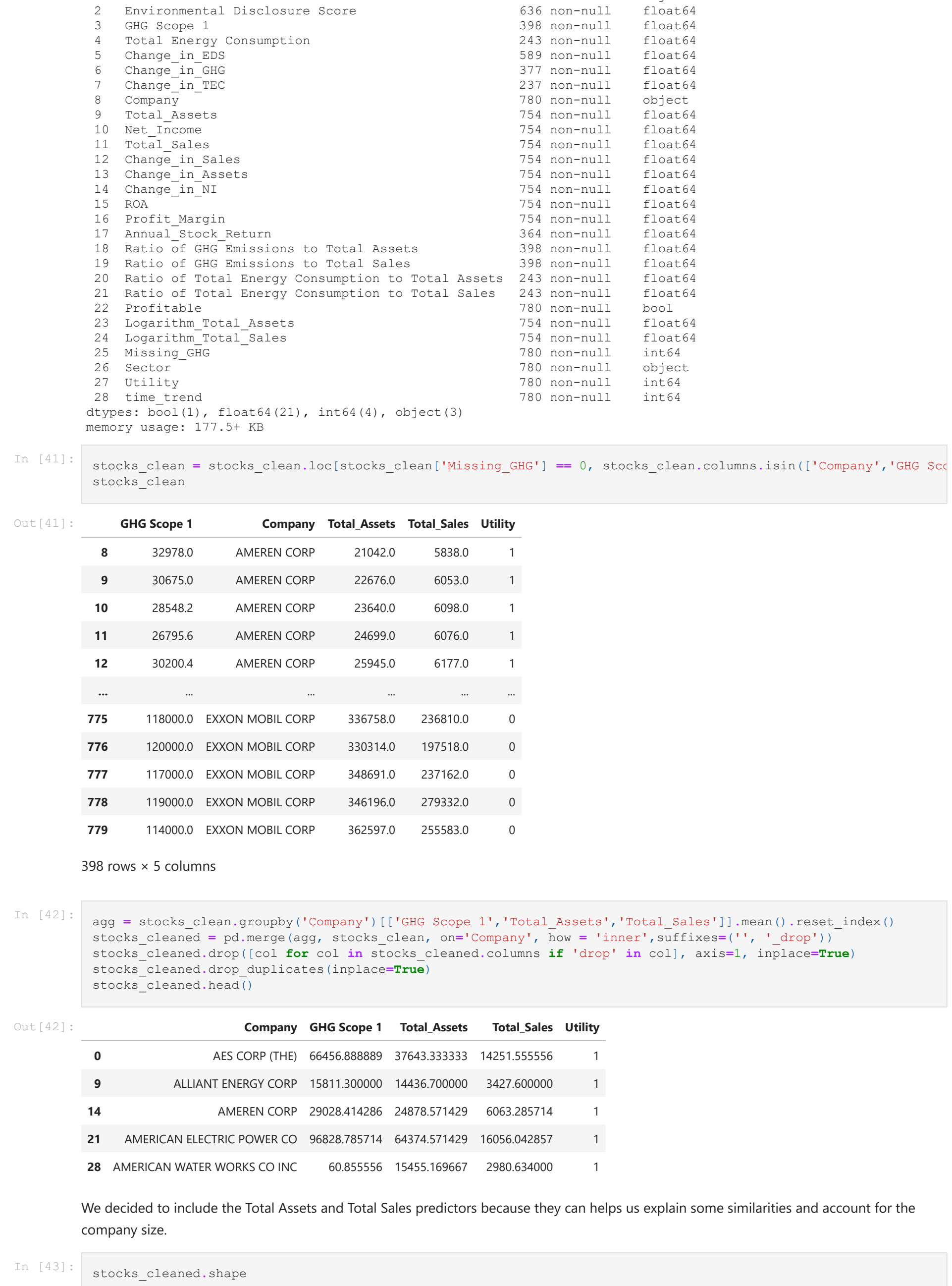
Cluster profiling for Energy industry



Data Cleaning - version 2

We decided to not consider the first data cleaning version for the cluster analysis. Our second and final version consists in just using the companies that are not missing values for the GHG Scope 1. We decided to do this to be consistent with the predictive models that were built in Sprint.k4.As a reminder, in that sprint, we build multiple models to predict the GHG Scope 1 and we only used the observations that didn't have null values.

After subsetting for the companies that have GHG Scope data, we collected the unique tickers and calculate the average for the features that we are insured to collect. For this version, we are going to consider the features GHG Scope 1, Total Assets, and Total Sales. Our goal is to be able to find patterns, similarities, and groups for these stocks.



Hierarchical Clustering - Version 2 and utility industry

We follow the same steps as with version 1.

1) Standardized data 2) Clustering using euclidean and cosine for distance matrix 3) Try different linkage method for each distance matrix 4) Create labels and add them to the original dataset 5) Sort by cluster values



K-Means Clustering - Version 2 and utility industry

#Cluster Evaluation - Deciding how many clusters

X = stock_scaled.values

KRANGE = range(2,10)

sse = []

Loop over and evaluate

for k in KRANGE:

km = KMeans(k)

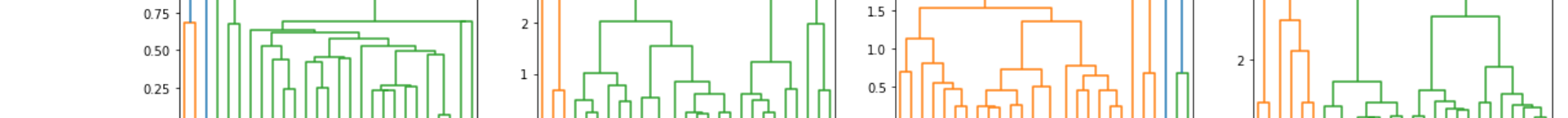
lab = km.fit_predict(stock_scaled)

sse.append(km.inertia_)

#Elbow Method

sns.lineplot(KRANGE,sse)

<AxesSubplot>



Testing K

ssl = []

for k in KRANGE:

km = KMeans(k)

lab = km.fit_predict(stock_scaled)

ssl.append(metrics.silhouette_score(stock_scaled, lab))

sns.lineplot(KRANGE, ssl)

<AxesSubplot>



K Means for 3

k3 = KMeans(3)

k3.fit(X)

k3_labels = k3.predict(X)

stocks_clean1['k3'] = k3_labels

stocks_clean1.k3.value_counts(dropna=False, sort=False) #

0 5

1 9

2 13

Name: k3, dtype: int64

We will proceed to export this utility dataset which will be used later on this notebook and in other ones.

stocks_clean1.to_csv('/Users/marinaltorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/DataSet')

Hierarchical Clustering - Version 2 and energy industry

We follow the same steps as with version 2 and utility industry.

1) Standardized data 2) Clustering using euclidean and cosine for distance matrix 3) Try different linkage method for each distance matrix 4) Create labels and add them to the original dataset 5) Sort by cluster values



K-means Clustering - Version 2 and energy industry

#Cluster Evaluation - Deciding how many clusters

X = stock_scaled.values

KRANGE = range(2,10)

sse = []

Loop over and evaluate

for k in KRANGE:

km = KMeans(k)

lab = km.fit_predict(stock_scaled)

sse.append(km.inertia_)

#Elbow Method

sns.lineplot(KRANGE,sse)

<AxesSubplot>



Testing K

ssl = []

for k in KRANGE:

km = KMeans(k)

lab = km.fit_predict(stock_scaled)

ssl.append(metrics.silhouette_score(stock_scaled, lab))

sns.lineplot(KRANGE, ssl)

<AxesSubplot>



k3 = KMeans(3)

k3.fit(X)

k3_labels = k3.predict(X)

stocks_clean2['k3'] = k3_labels

stocks_clean2.k3.value_counts(dropna=False, sort=False) #

1 16

2 2

3 6

Name: k3, dtype: int64

Not good at all... We will proceed to export this utility dataset which will be used later on this notebook and in other ones.

stocks_clean2.to_csv('/Users/marinaltorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/DataSet')

Cluster profiling for version 2 and utility industry

df_clus = stocks_clean1.drop(columns=['k3','Utility'])

stock_numeric = df_clus.select_dtypes('number')

stock_numeric

