# Unsupervised Learning Project: Trade & Ahead

## Context

There are a lot of compelling reasons to invest in stocks. It can help in fighting inflation, create wealth, and also provides some tax benefits. Good steady returns on investments over a long period of time can also grow a lot more than seems possible. Also, thanks to the power of compound interest, the earlier one starts investing, the larger the corpus one can have for retirement. Overall, investing in stocks can help meet life's financial aspirations.

It is important to maintain a diversified portfolio when investing in stocks in order to maximise earnings under any market condition. Having a diversified portfolio tends to yield higher returns and face lower risk by tempering potential losses when the market is down. It is often easy to get lost in a sea of financial metrics to analyze while determining the worth of a stock, and doing the same for a multitude of stocks to identify the right picks for an individual can be a tedious task. By doing a cluster analysis, one can identify stocks that exhibit similar characteristics and ones which exhibit minimum correlation. This will help investors better analyze stocks across different market segments and help protect against risks that could make the portfolio vulnerable to losses.

## Objective

Trade&Ahead is a financial consultancy firm who provide their customers with personalized investment strategies. They have hired you as a Data Scientist and provided you with data comprising stock price and some financial indicators for a few companies listed under the New York Stock Exchange. They have assigned you the tasks of analyzing the data, grouping the stocks based on the attributes provided, and sharing insights about the characteristics of each group.

# **Data Dictionary**

- Ticker Symbol: An abbreviation used to uniquely identify publicly traded shares of a particular stock on a particular stock market
- Company: Name of the company
- GICS Sector: The specific economic sector assigned to a company by the Global Industry Classification Standard (GICS) that best defines its business operations
- GICS Sub Industry: The specific sub-industry group assigned to a company by the Global Industry Classification Standard (GICS) that best defines its business operations

- Current Price: Current stock price in dollars
- Price Change: Percentage change in the stock price in 13 weeks
- Volatility: Standard deviation of the stock price over the past 13 weeks
- ROE: A measure of financial performance calculated by dividing net income by shareholders' equity (shareholders' equity is equal to a company's assets minus its debt)
- Cash Ratio: The ratio of a company's total reserves of cash and cash equivalents to its total current liabilities
- Net Cash Flow: The difference between a company's cash inflows and outflows (in dollars)
- Net Income: Revenues minus expenses, interest, and taxes (in dollars)
- Earnings Per Share: Company's net profit divided by the number of common shares it has outstanding (in dollars)
- Estimated Shares Outstanding: Company's stock currently held by all its shareholders
- P/E Ratio: Ratio of the company's current stock price to the earnings per share
- P/B Ratio: Ratio of the company's stock price per share by its book value per share (book value of a company is the net difference between that company's total assets and total liabilities)

## Importing necessary libraries, data, and creating functions

```
1 from google.colab import drive
 2 drive.mount('/content/drive')
→ Mounted at /content/drive
 1 # Libraries to help with reading and manipulating data
 2 import numpy as np
 3 import pandas as pd
 5 # Libraries to help with data visualization
 6 import matplotlib.pyplot as plt
 7 import seaborn as sns
 9 # to scale the data using z-score
10 from sklearn.preprocessing import StandardScaler
11
12 # to compute distances
13 from scipy.spatial.distance import cdist, pdist
14
15 # to perform k-means clustering and compute silhouette scores
16 from sklearn.cluster import KMeans
17 from sklearn.metrics import silhouette_score
18
19 # to perform hierarchical clustering, compute cophenetic correlation, and crea
```

```
20 from sklearn.cluster import AgglomerativeClustering
21 from scipy.cluster.hierarchy import dendrogram, linkage, cophenet
22
23 # to perform PCA
24 from sklearn.decomposition import PCA
26 # to visualize the elbow curve and silhouette scores
27 from yellowbrick.cluster import KElbowVisualizer, SilhouetteVisualizer
29 #format numeric data for easier readability
30 pd.set_option(
      "display.float_format", lambda x: "%.2f" % x
32 ) # to display numbers rounded off to 2 decimal places
33
34 # Removes the limit for the number of displayed columns
35 pd.set_option("display.max_columns", None)
36 # Sets the limit for the number of displayed rows
37 pd.set_option("display.max_rows", 200)
38
39 # To supress warnings
40 import warnings
41 warnings.filterwarnings("ignore")
```

```
1 # function to create labeled barplots
 2
 3
 4 def labeled barplot(data, feature, perc=False, n=None):
 5
 6
       Barplot with percentage at the top
 7
       data: dataframe
 8
       feature: dataframe column
 9
       perc: whether to display percentages instead of count (default is False)
10
       n: displays the top n category levels (default is None, i.e., display all levels)
11
12
13
14
       total = len(data[feature]) # length of the column
       count = data[feature].nunique()
15
16
       if n is None:
17
           plt.figure(figsize=(count + 1, 5))
       else:
18
19
           plt.figure(figsize=(n + 1, 5))
20
21
       plt.xticks(rotation=90, fontsize=12)
22
       ax = sns.countplot(
           data=data,
23
24
           x=feature,
25
           palette="viridis",
26
           order=data[feature].value_counts().index[:n].sort_values(),
27
       )
28
29
       for p in ax.patches:
30
           if perc == True:
               label = "{:.1f}%".format(
31
32
                   100 * p.get height() / total
               ) # percentage of each class of the category
33
34
           else:
35
               label = p.get_height() # count of each level of the category
36
37
           x = p.get_x() + p.get_width() / 2 # width of the plot
           y = p.get_height() # height of the plot
38
39
40
           ax.annotate(
41
               label,
               (x, y),
42
43
               ha="center",
44
               va="center",
45
               size=12,
46
               xytext=(0, 5),
47
               textcoords="offset points",
48
           ) # annotate the percentage
49
       plt.show() # show the plot
50
```

```
1 # function to plot a boxplot and a histogram along the same scale
 3 def histogram_boxplot(data, feature, figsize=(16, 6), kde=False, bins=None, hue=None):
 4
 5
       Combines boxplot and histogram
 6
 7
       data: dataframe
       feature: dataframe column
 8
       figsize: size of figure (default (16,6))
 9
       kde: whether to show the density curve (default False)
10
       bins: number of bins for histogram (default None)
11
12
       f2, (ax_box2, ax_hist2) = plt.subplots(
13
           nrows=2, # Number of rows of the subplot grid= 2
14
15
           sharex=True, # x-axis will be shared among all subplots
           gridspec_kw={"height_ratios": (0.25, 0.75)},
16
17
           figsize=figsize,
       ) # creating the 2 subplots
18
19
       sns.boxplot(
20
           data=data, x=feature, ax=ax_box2, showmeans=True,
21
       ) # boxplot will be created and a star will indicate the mean value of the column
22
       sns.histplot(
           data=data, x=feature, kde=kde, ax=ax_hist2, bins=bins, palette="winter",
23
24
       ) if bins else sns.histplot(
25
           data=data, x=feature, kde=kde, ax=ax_hist2
       ) # For histogram
26
27
       ax_hist2.axvline(
           data[feature].mean(), color="green", linestyle="--"
28
       ) # Add mean to the histogram
29
30
       ax hist2.axvline(
           data[feature].median(), color="black", linestyle="-"
31
       ) # Add median to the histogram
32
 1 #reading csv file into a pandas Dataframe
 2 ta = pd.read csv('/content/drive/MyDrive/Unsupervised Learning/stock data.csv')
 3 # copying data to another varaible to preserve original data
 4 df = ta.copy()
```

## Data Overview

```
1 # print a sample of five rows randomly selected from the training data
2 df.sample(n=5)
```



	Ticker Symbol	Security	GICS Sector	GICS Sub Industry	Current Price	Price Change	Volatility	ROE	Cash Ratio
75	COF	Capital One Financial	Financials	Consumer Finance	72.18	-0.62	1.36	9	99
199	MAS	Masco Corp.	Industrials	Building Products	28.30	11.64	1.43	263	61
214	MPC	Marathon Petroleum	Energy	Oil & Gas Refining & Marketing & Transportation	51.84	11.51	1.99	22	18
235	OKE	ONEOK	Energy	Oil & Gas Exploration & Production	24.66	-24.12	3.56	73	6
334	XYL	Xylem Inc.	Industrials	Industrial Conglomerates	36.50	11.01	1.17	16	83
4									•

1 df.shape

**→** (340, 15)

• Dataset has 340 rows and 15 columns

1 # print the data types of the columns within the datset
2 df.info()

<<class 'pandas.core.frame.DataFrame'>
 RangeIndex: 340 entries, 0 to 339
 Data columns (total 15 columns):

#	Column	Non-Null Count	Dtype
0	Ticker Symbol	340 non-null	object
1	Security	340 non-null	object
2	GICS Sector	340 non-null	object
3	GICS Sub Industry	340 non-null	object
4	Current Price	340 non-null	float64
5	Price Change	340 non-null	float64
6	Volatility	340 non-null	float64
7	ROE	340 non-null	int64
8	Cash Ratio	340 non-null	int64
9	Net Cash Flow	340 non-null	int64
10	Net Income	340 non-null	int64
11	Earnings Per Share	340 non-null	float64
12	Estimated Shares Outstanding	340 non-null	float64
13	P/E Ratio	340 non-null	float64

```
14 P/B Ratio 340 non-null float64 dtypes: float64(7), int64(4), object(4) memory usage: 40.0+ KB

1 # checking for duplicate values 2 df.duplicated().sum()

→ 0
```

- Dataset has no missing or duplicate values
- All columns with dtype object should be dtype category in order to conserve memory

```
1 # convert all columns with dtype object into category
2 for col in df.columns[df.dtypes=='object']:
      df[col] = df[col].astype('category')
1 # dropping the ticker symbol column, as it does not provide any information
2 df.drop("Ticker Symbol", axis=1, inplace=True)
1 # confirm new dataset
2 df.info()
→▼ <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 340 entries, 0 to 339
    Data columns (total 14 columns):
        Column
                                     Non-Null Count Dtype
    --- -----
                                     -----
                                     340 non-null
     0 Security
                                                    category
        GICS Sector
                                     340 non-null
                                                    category
                                     340 non-null
     2 GICS Sub Industry
                                                    category
                                     340 non-null float64
       Current Price
                                     340 non-null float64
        Price Change
     5
                                     340 non-null float64
        Volatility
                                    340 non-null int64
        ROE
     7
        Cash Ratio
                                    340 non-null
                                                    int64
       Net Cash Flow
                                                    int64
                                     340 non-null
        Net Income
                                     340 non-null
                                                    int64
     10 Earnings Per Share
                                     340 non-null float64
     11 Estimated Shares Outstanding 340 non-null
                                                    float64
     12 P/E Ratio
                                     340 non-null
                                                    float64
     13 P/B Ratio
                                     340 non-null
                                                    float64
```

- The 14 columns have three different dtypes: category(3), float64(7), int64(4)
- All of these dtypes are appropriate for their respective columns

dtypes: category(3), float64(7), int64(4)

memory usage: 46.7 KB

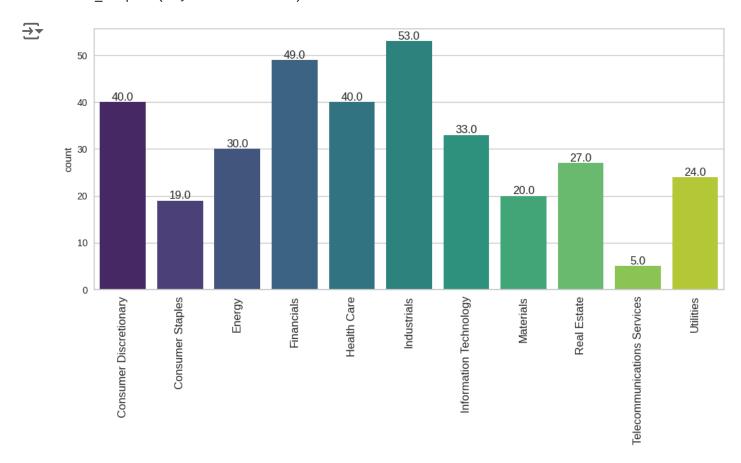
# Exploratory Data Analysis (EDA)

- 1 #provide statistical summary of all categorical columns
- 2 df.describe(include='category').T

<b>→</b>		count	unique	top	freq
	Security	340	340	3M Company	1
	GICS Sector	340	11	Industrials	53
	GICS Sub Industry	340	104	Oil & Gas Exploration & Production	16

#### **GICS Sector**

- 1 #create labeled barplot of stocks by sector
- 2 labeled\_barplot(df, 'GICS Sector')



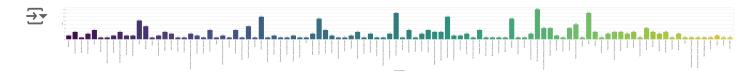
- 1 #display the five sectors with the most number of stocks
  2 df["GICS Sector"].value\_counts().head(n=5)
- Financials 53
  For Consumer Discretionary 40

Health Care 40
Information Technology 33
Name: GICS Sector, dtype: int64

- The stocks are drawn from 11 different industrial sectors, with no one sector comprising more than 16% of the dataset
- The top 4 of the 11 sectors (industrials, financials, consumer discretionary, and health care) comprise over half of the total number of stocks

## **GICS Sub Industry**

```
1 #create labeled barplot of stocks by sub industry
2 labeled_barplot(df, 'GICS Sub Industry')
```



1 #display the five sub industries with the most number of stocks
2 df['GICS Sub Industry'].value\_counts().head(n=5)

Oil & Gas Exploration & Production
REITS 14
Industrial Conglomerates 14
Internet Software & Services 12
Electric Utilities 12
Name: GICS Sub Industry, dtype: int64

- The dataset is comprised of stocks from 104 different subindustries, with no subindustry having more than 16 stocks in the dataset
- These observations indicate that the 340 stocks held within the dataset are highly diversified across sectors and subindustries

```
1 #provide statistical summary of all numerical columns
2 df.describe().T
```



	count	mean	std	min	25%	
Current Price	340.00	80.86	98.06	4.50	38.55	ţ
Price Change	340.00	4.08	12.01	-47.13	-0.94	
Volatility	340.00	1.53	0.59	0.73	1.13	
ROE	340.00	39.60	96.55	1.00	9.75	,
Cash Ratio	340.00	70.02	90.42	0.00	18.00	۷
Net Cash Flow	340.00	55537620.59	1946365312.18	-11208000000.00	-193906500.00	209800
Net Income	340.00	1494384602.94	3940150279.33	-23528000000.00	352301250.00	70733600
Earnings Per Share	340.00	2.78	6.59	-61.20	1.56	
Estimated Shares Outstanding	340.00	577028337.75	845849595.42	27672156.86	158848216.10	30967513
P/E Ratio	340.00	32.61	44.35	2.94	15.04	2
P/B Ratio	340.00	-1.72	13.97	-76.12	-4.35	
•						•

#### **Numerical Columns**

```
1 #create list of columns with numerical variables
```

3

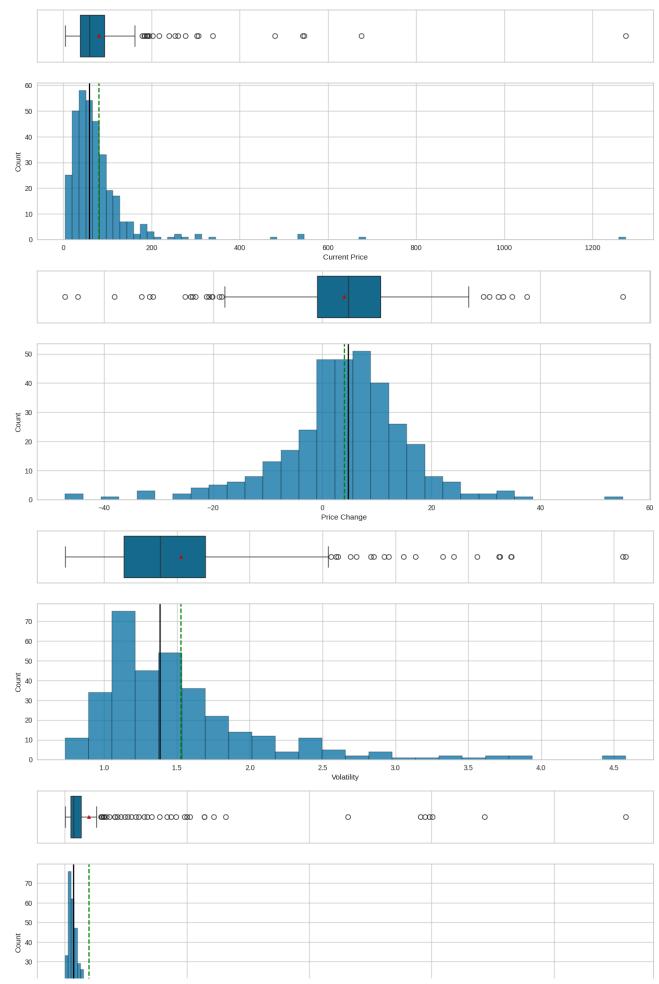
<sup>2</sup> num\_col = df.select\_dtypes(include=np.number).columns.tolist()

<sup>4 #</sup>display histograms and boxplots for all numerical columns

<sup>5</sup> for col in num\_col:

<sup>6</sup> histogram\_boxplot(df, col)





## **Current price**

- The distribution is heavily right skewed, with 49 of the 340 stocks having twice the median value of all stocks
- As expected, no stock is listed at less of less than 0 dollars

## **Price change**

- The distribution is biased towards lower volatilities, but long tails do exist both for positive and negative price changes
- The most volatile stocks show as low as a 47% decrease to as high as a 55% increase over 13 weeks

## **Volatility**

As expected, the distribution of standard deviations is right skewed and not normal

#### Cash Ratio / ROE

- As expected, both distributions are heavily right skewed and no stock is listed with either metric with a value of less than 0
  - For example, 24 stocks are listed with returns on equity of less than 5 and 25 stocks are listed with returns of over 100 percent

#### **Net Income / EPS**

- As expected, net income is shown to be right skewed with both long positive and negative tails
  - I.e., most companies generate meager profits, but some are failing and some are highly successful
- 32 companies within the dataset are showing a net income of less than 0 dollars
- EPS, as a derivative of Net Income, shows a similar distribution, with most showing low positive values and a few stocks (34) showing negative values

## **Estimated shares outstanding**

 The distribution is highly right skewed, but no stock has a value of outstanding shares that is unrealistic

#### P/E and P/B Ratio

- The distribution of P/E ratios is highly right skewed
  - Interestingly, no stock shows a negative ratio, even though several stocks have a negative EPS and no stock stock has a price listed of less than 0

- The distribution for P/B ratios is mostly centered around 0 but with long positive and negative
  - For example, 175 of the 340 total stocks are shown to below the 25th percentile and above the 75th percentile and
  - o Additionally, 31 of the stocks are outliers

## Conclusions

- As expected, stocks offer uncertain returns with high upsides, mostly modest returns, and the
  omnipresent possibility that the value of the stock may become worthless (i.e., the company
  goes bankrupt)
- All of these variables contain a few or several outliers; however, none of these values appear
  to be unrealistic given the nature of stock prices and historical expectations

## The stocks of which economic sector have seen the maximum price increase on average?

0 1 2 3 4 5 6

```
1 df.groupby('GICS Sector')['Price Change'].mean().sort_values()
```

$\rightarrow$	GICS Sector								
	Energy	-10.23							
	Utilities	0.80							
	Industrials	2.83							
	Financials	3.87							
	Materials	5.59							
	Consumer Discretionary	5.85							
	Real Estate	6.21							
	Telecommunications Services	6.96							
	Information Technology	7.22							
	Consumer Staples	8.68							
	Health Care	9.59							
	Name: Price Change, dtype: flo	oat64							

• Stocks within the health care sectors have shown the highest average price increase over the preeceding period



#### How are the different variables correlated with each other?

```
1 #create correlation heat map for numerical variables
 2 plt.figure(figsize=(14, 7))
 3 sns.heatmap(
 4
       df[num_col].corr(),
 5
       annot=True,
 6
       vmin=-1,
 7
       vmax=1,
       fmt=".2f",
 8
 9
       cmap='viridis'
10)
11 plt.show()
```





- Several variables are moderately correlated (+/- .40) with one another
  - Volatility is negatively correlated with price change, i.e., as a stock becomes more volatile, its price is likely dropping
  - Net income is negatively correlayed with volatility, i.e. as a company generates higher net income its price is likely less volatile
  - Net income is also positively correlated with earnings per share (EPS) and estimated shares outstanding
  - EPS is positively correlated with current price, i.e. as a company's EPS rises, its prices is also highly likely to increase
  - EPS is also negatively correlated with ROE, i.e. as a company generates more equity for shareholders, an equivalent amount of net income the following periods will generate a

lower return

Cash ratio provides a measure of a company's ability to cover its short-term obligations using only cash and cash equivalents. How does the average cash ratio vary across economic sectors?

```
1 df.groupby('GICS Sector')['Cash Ratio'].mean().sort values(ascending=False)
→ GICS Sector
    Information Technology
                                   149.82
    Telecommunications Services
                                   117.00
    Health Care
                                   103.78
    Financials
                                    98.59
    Consumer Staples
                                    70.95
    Energy
                                    51.13
    Real Estate
                                    50.11
                                    49.58
    Consumer Discretionary
    Materials
                                    41.70
    Industrials
                                    36.19
    Utilities
                                    13.62
    Name: Cash Ratio, dtype: float64
```

- IT and Telecommunications sectors, both relatively newer and unregulated industries, are able to generate significantly higher average cash ratios than their peer sectors
- · Utilities, a highly regulated industry, generates the lowest average cash ratios of all sectors

P/E ratios can help determine the relative value of a company's shares as they signify the amount of money an investor is willing to invest in a single share of a company per dollar of its earnings. How does the P/E ratio vary, on average, across economic sectors?

```
1 df.groupby('GICS Sector')['P/E Ratio'].mean().sort_values(ascending=False)
→ GICS Sector
                                   72.90
    Energy
    Information Technology
                                   43.78
    Real Estate
                                   43.07
    Health Care
                                   41.14
                                   35.21
    Consumer Discretionary
    Consumer Staples
                                   25.52
    Materials
                                   24.59
                                   18.72
    Utilities
    Industrials
                                   18.26
    Financials
                                   16.02
    Telecommunications Services
                                   12.22
    Name: P/E Ratio, dtype: float64
```

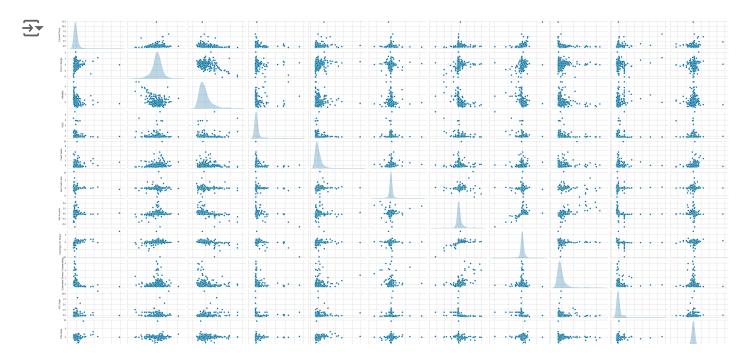
• Energy companies have the highest average P/E ratios of all sectors by a considerable margin, with telecoms having the lowest average P/E ratios

# K-means Clustering

```
1 #scale the data set before clustering
2 scaler = StandardScaler()
3 subset = df[num_col].copy()
4 subset_scaled = scaler.fit_transform(subset)

1 #create a dataframe from the scaled data
2 subset_scaled_df = pd.DataFrame(subset_scaled, columns=subset.columns)

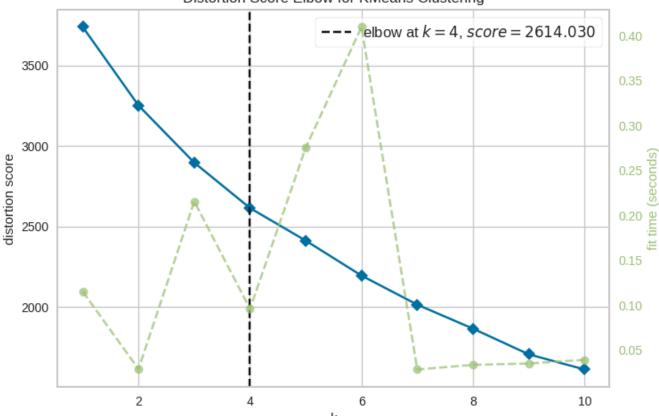
1 #create pairplot for scaled dataframe
2 sns.pairplot(subset_scaled_df, height=2,aspect=2 , diag_kind='kde')
3 plt.show()
```



```
1 #print average distortions for range of kmeans models fitted to scaled dataset
 2 clusters = range(1, 11)
 3 meanDistortions = []
 4
 5 for k in clusters:
 6
       model = KMeans(n clusters=k)
 7
       model.fit(subset_scaled_df)
       prediction = model.predict(subset_scaled_df)
 8
 9
       distortion = (
10
           sum(
11
               np.min(cdist(subset_scaled_df, model.cluster_centers_, "euclidean"), axis=1)
12
           / subset scaled df.shape[0]
13
       )
14
15
16
       meanDistortions.append(distortion)
17
18
       print("Number of Clusters:", k, "\tAverage Distortion:", distortion)
     Number of Clusters: 1
                             Average Distortion: 2.5425069919221697
     Number of Clusters: 2
                             Average Distortion: 2.382318498894466
     Number of Clusters: 3
                             Average Distortion: 2.2659465936501304
     Number of Clusters: 4
                             Average Distortion: 2.2106334639099248
     Number of Clusters: 5
                             Average Distortion: 2.128338996505974
     Number of Clusters: 6
                             Average Distortion: 2.0598763618226803
     Number of Clusters: 7
                             Average Distortion: 2.0246699618696655
     Number of Clusters: 8
                             Average Distortion: 1.9887592433275094
     Number of Clusters: 9
                             Average Distortion: 1.8929684909747697
     Number of Clusters: 10 Average Distortion: 1.8540958566681758
 1 #fit KMeans model and use visualizaer to indicate optimal K value
 2 model = KMeans(random state=42)
 3 visualizer = KElbowVisualizer(model, k=(1, 11), timings=True)
 4 visualizer.fit(subset scaled df) # fit the data to the visualizer
 5 visualizer.show() # finalize and render figure
 6 plt.show()
```

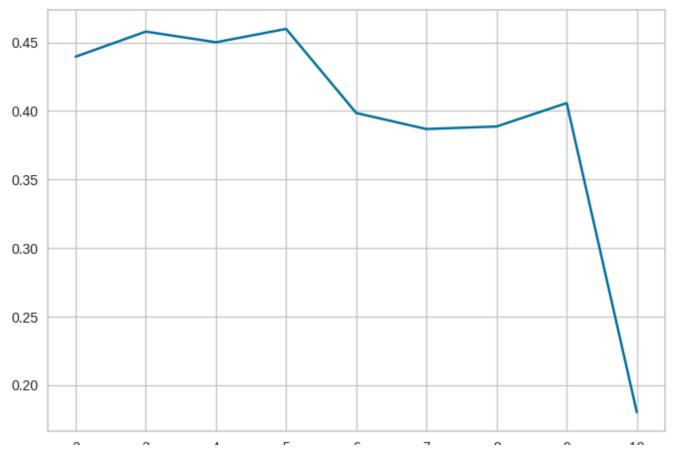
**₹** 

#### Distortion Score Elbow for KMeans Clustering



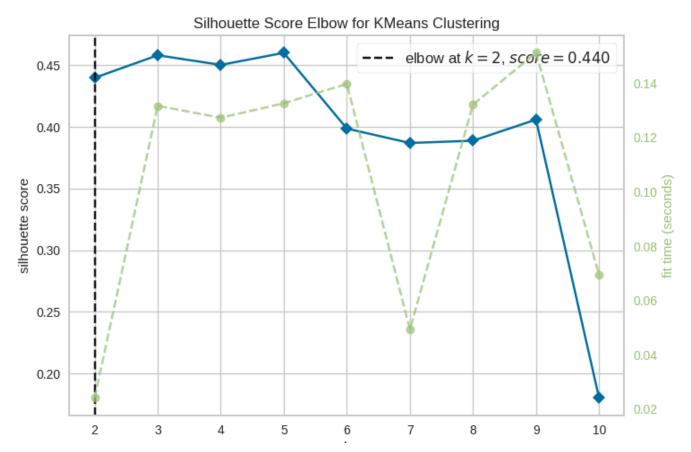
```
1 #fit KMeans model and provide silhouette scores for range of k clusters
 2 sil_score = []
 3 cluster_list = range(2, 11)
 4 for n_clusters in cluster_list:
       clusterer = KMeans(n_clusters=n_clusters, random_state=42)
 5
       preds = clusterer.fit_predict((subset_scaled_df))
 6
 7
       score = silhouette_score(subset_scaled_df, preds)
 8
       sil_score.append(score)
       print("For n_clusters = {}, the silhouette score is {})".format(n_clusters, score))
 9
10
11 #show scores in line graph
12 plt.plot(cluster_list, sil_score)
13 plt.show()
```

```
\rightarrow For n_clusters = 2, the silhouette score is 0.43969639509980457)
    For n_clusters = 3, the silhouette score is 0.45797710447228496)
    For n_clusters = 4, the silhouette score is 0.45017906939331087)
    For n_clusters = 5, the silhouette score is 0.4599352800740646)
    For n_clusters = 6, the silhouette score is 0.3985379248608659)
    For n_clusters = 7, the silhouette score is 0.3868475076242907)
    For n_clusters = 8, the silhouette score is 0.3886929719130642)
    For n_clusters = 9, the silhouette score is 0.40581042332267614)
    For n_clusters = 10, the silhouette score is 0.18011528994705786)
```



```
1 #fit KMeans model and use visualizaer to indicate optimal K value
2 model = KMeans(random_state=42)
3 visualizer = KElbowVisualizer(model, k=(2, 11), metric="silhouette", timings=True)
4 visualizer.fit(subset_scaled_df) # fit the data to the visualizer
5 visualizer.show() # finalize and render figure
6 plt.show()
```

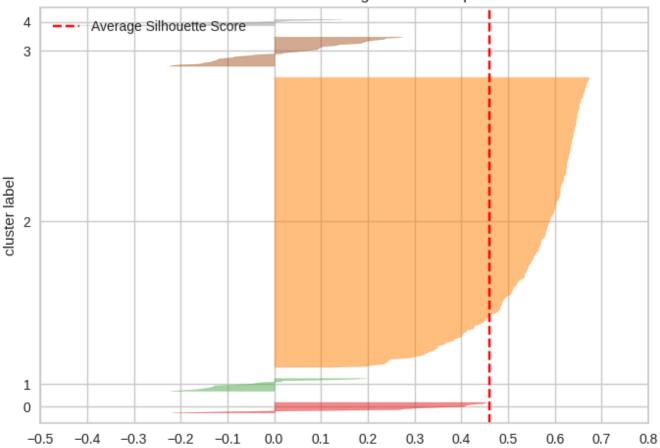




- 1 #find optimal no. of clusters with silhouette coefficients
- 2 visualizer = SilhouetteVisualizer(KMeans(5, random\_state=42))
- 3 visualizer.fit(subset\_scaled\_df)
- 4 visualizer.show()
- 5 plt.show()

 $\overline{2}$ 





 Between the Elbow and Silhouette plots, the number of clusters with the best performance appears to be 5

```
1 #create kmeans cluster model
2 kmeans = KMeans(n_clusters=5, random_state=42)
3
4 #fit model to scaled dataset
5 kmeans.fit(subset_scaled_df)

The scale of the scale o
```

## Cluster Profiling

```
1 # adding kmeans cluster labels to the original dataframe
2 df["KMeans_clusters"] = kmeans.labels_
```

```
1 #group dataset by kmeans cluster labels
2 cluster_profile = df.groupby("KMeans_clusters").mean()
3
4 #add counts for number of stocks in each cluster
5 cluster_profile["Count"] = (
6    df.groupby("KMeans_clusters")["Current Price"].count().values
7 )
1 cluster_profile.style.highlight_max(color="lightblue", axis=0)
```

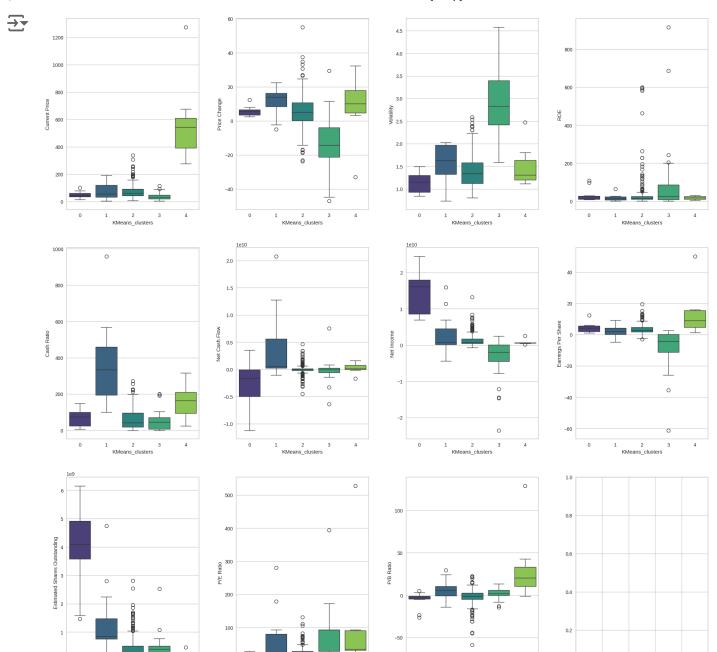
 $\overline{\mathbf{T}}$ 

	Current Price	Price Change	Volatility	ROE	Cash Ratio	Net Cash
KMeans_clusters						
0	50.477272	5.588148	1.141171	31.545455	64.181818	-2581727272.72
1	81.418719	10.536341	1.578634	17.000000	367.538462	3857062692.30
2	73.769121	5.466467	1.392827	34.632143	55.710714	4183132.14
3	38.808966	-13.680395	2.938240	106.034483	55.551724	-189825655.17
4	585.527134	7.752090	1.508020	17.571429	159.142857	210520428.57
4						•

```
1 # print the names of the companies in each cluster
2 for cl in df["KMeans_clusters"].unique():
      print("In cluster {}, the following companies are present:".format(cl))
      print(df[df["KMeans_clusters"] == cl]["Security"].unique().to_list())
4
5
      print()
→ In cluster 2, the following companies are present:
    ['American Airlines Group', 'AbbVie', 'Abbott Laboratories', 'Adobe Systems Inc', 'Arche
    In cluster 1, the following companies are present:
    ['Analog Devices, Inc.', 'Alexion Pharmaceuticals', 'Amgen Inc', 'Bank of America Corp',
    In cluster 4, the following companies are present:
    ['Alliance Data Systems', 'Amazon.com Inc', 'Chipotle Mexican Grill', 'Equinix', 'Intuit
    In cluster 3, the following companies are present:
    ['Apache Corporation', 'Anadarko Petroleum Corp', 'Baker Hughes Inc', 'Chesapeake Energy
    In cluster 0, the following companies are present:
    ['Citigroup Inc.', 'Ford Motor', 'Gilead Sciences', 'General Motors', 'JPMorgan Chase &
```

```
1 #print number of stocks within each sector for all of the clusters
2 for k in range(0,df['KMeans_clusters'].nunique()):
      print('The number of stocks within each GICS Sector for Cluster '+str(k)+' are:')
      print(df[df['KMeans clusters']==k]['GICS Sector'].value counts())
4
5
      print(" ")
\rightarrow \overline{\phantom{m}} The number of stocks within each GICS Sector for Cluster 0 are:
    Financials
                                     3
    Consumer Discretionary
                                     2
    Health Care
                                     2
    Telecommunications Services
                                     2
    Consumer Staples
    Energy
                                     1
    Industrials
    Information Technology
                                     0
    Materials
    Real Estate
    Utilities
    Name: GICS Sector, dtype: int64
    The number of stocks within each GICS Sector for Cluster 1 are:
    Health Care
    Information Technology
                                     4
    Consumer Discretionary
                                     1
    Consumer Staples
    Energy
                                     1
    Financials
    Telecommunications Services
    Industrials
    Materials
                                     0
    Real Estate
                                     0
    Utilities
    Name: GICS Sector, dtype: int64
    The number of stocks within each GICS Sector for Cluster 2 are:
    Industrials
                                     52
    Financials
                                     45
    Consumer Discretionary
                                     33
    Health Care
                                     32
    Real Estate
                                     26
    Information Technology
                                     25
    Utilities
                                     24
    Materials
                                     18
    Consumer Staples
                                     17
    Energy
                                      6
    Telecommunications Services
    Name: GICS Sector, dtype: int64
    The number of stocks within each GICS Sector for Cluster 3 are:
    Energy
                                     22
    Information Technology
                                      3
                                      2
    Materials
    Consumer Discretionary
                                      1
                                      1
    Industrials
    Consumer Staples
    Financials
```

```
Health Care
                                      0
     Real Estate
                                      0
     Telecommunications Services
                                      0
     Utilities
     Name: GICS Sector, dtype: int64
     The number of stocks within each GICS Sector for Cluster 4 are:
     Consumer Discretionary
                                     3
 1 # show boxplots of numerical variables for each K-Means cluster
 2 fig, axes = plt.subplots(3, 4, figsize=(20, 20))
 3 \text{ counter} = 0
 5 for ii in range(3):
 6
       for jj in range(4):
 7
           if counter < 11:
               sns.boxplot(
 8
 9
                   ax=axes[ii][jj],
                   data=df,
10
                   y=df.columns[3+counter],
11
12
                   x="KMeans_clusters",
13
                   palette="viridis"
14
15
               counter = counter + 1
16
17 fig.tight_layout(pad=3.0)
```



## **KMeans Clusters**

## Cluster 0 - Large Market Capitalization / Dow Jones Industrial Average

- 11 stocks, comprised mostly of stocks within the Financials, Health Care, Information Technology (IT), and Consumer Discretionary sectors
- · Companies within this cluster have:
  - Low volatility
  - o Most of the companies with the highest outflows of cash
  - The highest net incomes
  - The highest number of shares outstanding

#### Cluster 1 - "Cash is King"

- 13 stocks, comprised mostly of stocks within the Healthcare and IT sectors
- · Companies within this cluster have:
  - Moderate volatility
  - Mostly profitable
  - Most of the highest cash ratios and cash inflows

#### Cluster 2 - S&P 500 / Diversification

- 280 stocks (~84% of all stocks in the dataset) drawn from all sectors present in the dataset
- Companies within this cluster have:
  - Low P/E ratios
  - Most of the outliers on negative P/B ratios

#### Cluster 3 - "Ride the Energy Rollercoaster" portfolio / Growth mindset

- 29 stocks, a vast majority of which are from the Energy sector
- Companies within this cluster have:
  - Low stock prices, but high ROE
  - High beta
  - Most of the most volatile stocks, especially those with outliers in price decreases
  - Mostly negative net incomes and earnings per share

#### **Cluster 4 - High Earnings for a High Price**

- 7 stocks, comprised mostly of stocks from the Health Care and Consumer Discretionary sectors
- Companies within this cluster have:
  - Most of stocks with the highest prices
  - Favorable cash ratios
  - The most favorable P/B ratios
  - Most of the highest earnings-per-share

# Hierarchical Clustering

```
1 # list of distance metrics
 2 distance_metrics = ["euclidean", "chebyshev", "mahalanobis", "cityblock"]
 4 # list of linkage methods
 5 linkage_methods = ["single", "complete", "average", "weighted"]
 7 high_cophenet_corr = 0
 8 \text{ high\_dm\_lm} = [0, 0]
 9
10 for dm in distance_metrics:
       for lm in linkage_methods:
11
           Z = linkage(subset_scaled_df, metric=dm, method=lm)
12
           c, coph_dists = cophenet(Z, pdist(subset_scaled_df))
13
14
           print(
15
               "Cophenetic correlation for {} distance and {} linkage is {}.".format(
                   dm.capitalize(), lm, round(c,4)
16
17
               )
18
           print(" ")
19
20
           if high_cophenet_corr < c:</pre>
               high_cophenet_corr = c
21
               high_dm_lm[0] = dm
22
23
               high_dm_lm[1] = lm
```

 $\rightarrow$ Cophenetic correlation for Euclidean distance and single linkage is 0.9232. Cophenetic correlation for Euclidean distance and complete linkage is 0.7873. Cophenetic correlation for Euclidean distance and average linkage is 0.9423. Cophenetic correlation for Euclidean distance and weighted linkage is 0.8694. Cophenetic correlation for Chebyshev distance and single linkage is 0.9063. Cophenetic correlation for Chebyshev distance and complete linkage is 0.5989. Cophenetic correlation for Chebyshev distance and average linkage is 0.9338. Cophenetic correlation for Chebyshev distance and weighted linkage is 0.9127. Cophenetic correlation for Mahalanobis distance and single linkage is 0.9259. Cophenetic correlation for Mahalanobis distance and complete linkage is 0.7925. Cophenetic correlation for Mahalanobis distance and average linkage is 0.9247. Cophenetic correlation for Mahalanobis distance and weighted linkage is 0.8708. Cophenetic correlation for Cityblock distance and single linkage is 0.9334. Cophenetic correlation for Cityblock distance and complete linkage is 0.7375. Cophenetic correlation for Cityblock distance and average linkage is 0.9302.

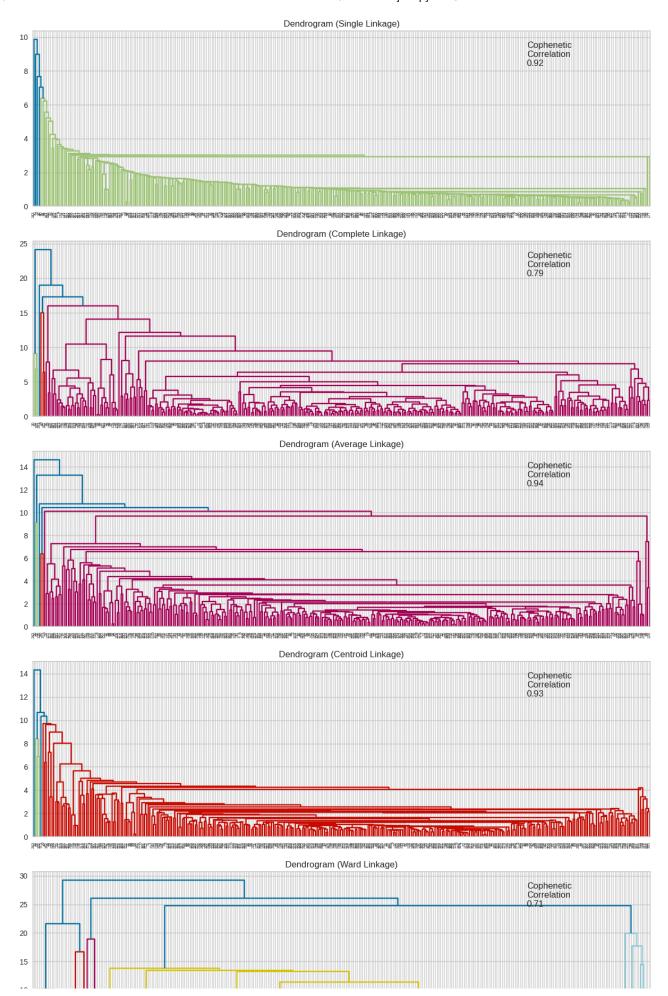
Cophenetic correlation for Cityblock distance and weighted linkage is 0.731.

1 # printing the combination of distance metric and linkage method with the highest cophene

```
2 print(
 3
       "Highest cophenetic correlation is {}, which is obtained with {} distance and {} link
 4
           round(high_cophenet_corr,4), high_dm_lm[0].capitalize(), high_dm_lm[1]
 5
       )
 6)
→ Highest cophenetic correlation is 0.9423, which is obtained with Euclidean distance and
 1 # list of linkage methods for euclidean distance metric
 2 linkage_methods = ["single", "complete", "average", "centroid", "ward", "weighted"]
 4 high_cophenet_corr = 0
 5 \text{ high\_dm\_lm} = [0, 0]
 7 for lm in linkage methods:
       Z = linkage(subset_scaled_df, metric="euclidean", method=lm)
       c, coph_dists = cophenet(Z, pdist(subset_scaled_df))
 9
       print(
10
               "Cophenetic correlation for Euclidean distance and {} linkage is {}.".format(
11
12
                   lm, round(c,4)
13
14
           )
       print(" ")
15
       if high_cophenet_corr < c:</pre>
16
17
           high cophenet corr = c
18
           high_dm_lm[0] = "euclidean"
19
           high dm lm[1] = lm
     Cophenetic correlation for Euclidean distance and single linkage is 0.9232.
\rightarrow
     Cophenetic correlation for Euclidean distance and complete linkage is 0.7873.
     Cophenetic correlation for Euclidean distance and average linkage is 0.9423.
     Cophenetic correlation for Euclidean distance and centroid linkage is 0.9314.
     Cophenetic correlation for Euclidean distance and ward linkage is 0.7101.
     Cophenetic correlation for Euclidean distance and weighted linkage is 0.8694.
```

```
1 # printing the combination of distance metric and linkage method with the highest cophene
 2 print(
       "Highest cophenetic correlation is {}, which is obtained with {} linkage.".format(
 3
 4
           round(high_cophenet_corr,4), high_dm_lm[1]
       )
 6)
\rightarrow \overline{\phantom{a}} Highest cophenetic correlation is 0.9423, which is obtained with average linkage.
 1 # list of linkage methods
 2 linkage_methods = ["single", "complete", "average", "centroid", "ward", "weighted"]
 4 # lists to save results of cophenetic correlation calculation
 5 compare_cols = ["Linkage", "Cophenetic Coefficient"]
 7 # to create a subplot image
 8 fig, axs = plt.subplots(len(linkage_methods), 1, figsize=(15, 30))
10 # We will enumerate through the list of linkage methods above
11 # For each linkage method, we will plot the dendrogram and calculate the cophenetic corre
12 for i, method in enumerate(linkage_methods):
       Z = linkage(subset scaled df, metric="euclidean", method=method)
13
14
15
       dendrogram(Z, ax=axs[i])
       axs[i].set_title(f"Dendrogram ({method.capitalize()} Linkage)")
16
17
       coph corr, coph dist = cophenet(Z, pdist(subset scaled df))
18
19
       axs[i].annotate(
20
           f"Cophenetic\nCorrelation\n{coph_corr:0.2f}",
21
           (0.80, 0.80),
           xycoords="axes fraction",
22
23
       )
```







- The cophenetic correlation is highest for average and centroid linkage methods, but the dendrogram for average appears to provide better clusters
- 5 appears to be the appropriate number of clusters for the average linkage method

## Cluster Profiling

```
1 df_hierarchy = df.copy()
2 df_hierarchy.drop("KMeans_clusters", axis=1, inplace=True)
3 df_hierarchy['HC_clusters'] = hierarchy.labels_

1 #group dataset by Hierarchical clusters
2 cluster_profile_h = df_hierarchy.groupby("HC_clusters").mean()
3
4 #add counts for number of stocks in each cluster
5 cluster_profile_h["Count"] = (
6     df_hierarchy.groupby("HC_clusters")["Current Price"].count().values
7 )
8
9 #show dataframe with maximum values for each metric highlighted
10 cluster_profile_h.style.highlight_max(color="lightblue", axis=0)
```



	Current Price	Price Change	Volatility	ROE	Cash Ratio	Net Cash Fl
HC_clusters						
0	77.884243	4.105986	1.516865	35.320359	66.775449	-32825817.3652
1	25.640000	11.237908	1.322355	12.500000	130.500000	16755500000.0000
2	24.485001	-13.351992	3.482611	802.000000	51.000000	-1292500000.0000
3	104.660004	16.224320	1.320606	8.000000	958.000000	592000000.0000
4	1274.949951	3.190527	1.268340	29.000000	184.000000	-1671386000.0000
4						•

- There are 2 clusters of one company, 2 clusters of two companies, and a single cluster of the remaining 334 companies
- The clustering of these companies does not solve the business problem at hand, because the clusters do not have enough variability

In contrasts, the dendrogram for Ward linkage appears to provide better clustering, with 5 appearing to be the appropriate number of clusters

# Cluster Profiling



	Current Price	Price Change	Volatility	ROE	Cash Ratio	Net Cash Flow
HC_clu	sters					
0	326.198218	10.563242	1.642560	14.400000	309.466667	288850666.666667
1	84.355716	3.854981	1.827670	633.571429	33.571429	-568400000.000000
2	42.848182	6.270446	1.123547	22.727273	71.454545	558636363.636364
3	72.760400	5.213307	1.427078	25.603509	60.392982	79951512.280702
4	36.440455	-16.073408	2.832884	57.500000	42.409091	-472834090.909091
4						•

```
1 # print the names of the companies in each cluster
2 for cl in df_hierarchy["HC_clusters"].unique():
      print("In cluster {}, the following companies are present:".format(cl))
      print(df_hierarchy[df_hierarchy["HC_clusters"] == cl]["Security"].unique().to_list())
5
      print()
\rightarrow In cluster 3, the following companies are present:
    ['American Airlines Group', 'AbbVie', 'Abbott Laboratories', 'Adobe Systems Inc', 'Analc
    In cluster 0, the following companies are present:
    ['Alliance Data Systems', 'Alexion Pharmaceuticals', 'Amgen Inc', 'Amazon.com Inc', 'Chi
    In cluster 1, the following companies are present:
    ['Allegion', 'Apache Corporation', 'Chesapeake Energy', 'Charter Communications', 'Colga
    In cluster 4, the following companies are present:
    ['Anadarko Petroleum Corp', 'Baker Hughes Inc', 'Cabot Oil & Gas', 'Concho Resources', '
    In cluster 2, the following companies are present:
    ['Bank of America Corp', 'Citigroup Inc.', 'Ford Motor', 'Intel Corp.', 'JPMorgan Chase
```

1 # print the number of stocks in each GICS sector for each cluster
2 for k in range(0,df\_hierarchy['HC\_clusters'].nunique()):
3 print('The number of stocks within each GICS Sector for Cluster '+str(k)+' are:')
4 print(df\_hierarchy[df\_hierarchy['HC\_clusters']==k]['GICS Sector'].value\_counts())
5 print(" ")

The number of stocks within each GICS Sector for Cluster 0 are:
Health Care 5
Information Technology 4
Consumer Discretionary 3
Consumer Staples 1
Real Estate 1
Telecommunications Services 1

Energy

7/10/24, 12:16 PM

Financials

0