

In [ ]:

```
1  #####
2  #      Project 4 :Assignment on Project
3  #
4  #####
5  # Loading the Data from the CSV file
6
7  # importing the packages and Loading the Data
8  import numpy as np
9  import pandas as pd
10 import matplotlib.pyplot as plt
11 % matplotlib inline
12
13 data = pd.read_csv("C:/Users/Maximus/pythonsessions/session31/data_stocks.csv")
14
15 print(data.head(2))
```

In [ ]:

```
1  # Analysing the Stock Data , Checking for Null values , variation , missing values
2
3  print( " Analysing the Data and knowing about the data")
4  print('-'*80)
5
6  print(data.info())
7
8  print('-'*80)
9
10 print(data.describe())
11 print('-'*80)
12
13 print(data.shape)
14 print('-'*80)
15
16 print(data.isnull().sum())
17 print('-'*80)
18 print( " We observe there are no null values \n", '-'*80)
```

```
In [ ]: 1 # As there are many columns , to get appropriate features , we do the feature selection first
2 # and removing the unwanted Features
3
4 # Since Date and SP500 columns have high std as compared to others as seen above, we can remove the columns
5 # Date is in epoch format
6
7 data.drop(['DATE', 'SP500'],axis=1,inplace=True)
8
```

**problem 1 : To Show all stocks are apparently similar in performance**

**Solution : So here we do Factor Analysis to find the similar stocks**

```
In [ ]: 1 # Assumption: As the Data Size huge a Sampling of 1000 records is considered
2
3 # We use Factor analysis here to get similar columns
4
5
6 from sklearn.decomposition import FactorAnalysis
7 from sklearn.model_selection import cross_val_score
8
9 # numFac number of Factors
10 n_components = np.arange(2,200)
11 fa_scores = []
12
13 # to find the number of factors
14 for numFac in n_components:
15     # define the Constructor for Factor Analysis
16     fa = FactorAnalysis(random_state=101)
17     fa.n_components = numFac
18     acc_score = np.mean(cross_val_score(fa,data.sample(1000), cv=5))
19     fa_scores.append((numFac,acc_score))
20
21
```

```
In [ ]: 1 # splitting the tuples into separate values of components and scores
2
3 n_compon, fa_score = zip(*fa_scores)
4
5 bstscore = fa_score[np.argmax(fa_score)]
6 bst_component = n_compon[np.argmax(fa_score)]
7 print(" The best Value of the score for the components is {} at components  {}\n ".format(bstscore, bst_component))
8 print('-'*80)
9
```

```
In [ ]: 1 # Visualising the scores and component value by plotting the graph for fa_Scores
2
3 # ncompon = n_compon[25:]
4 # fa_score1 = fa_scores[25:]
5
6 title = " N_components Vs Cv Score "
7 plt.figure()
8 plt.plot(n_compon[25:], fa_score[25:], 'r', label='FA scores')
9 plt.axvline(bst_component, color='r',
10             label='FactorAnalysis CV: %d' %bst_component,
11             linestyle='--')
12 plt.xlabel('nb of components')
13 plt.ylabel('CV scores')
14 plt.legend(loc='lower right')
15 plt.title(title)
16
17 plt.show()
18
```

```
In [ ]: 1 # To find the factors at the obtained best components value achieved
2 # Approach 1
3
4
5 factor = FactorAnalysis(n_components=bst_component , random_state=101).fit(data.sample(1000))
6 #print(factor.components_)
7
8 print('' Following is the DataFrame showing the relationship of columns Stock,
9       Those with positive values are similiar \n'', '-'*80)
10 fact_data = pd.DataFrame(factor.components_, columns=data.columns)
11 print(fact_data.head(5))
12
```

```
In [ ]: 1 # Another approach to show in each Factor the stocks that are similiar
2 from factor_analyzer import FactorAnalyzer
3 fa = FactorAnalyzer()
4
5 fa.analyze(data, bst_component, rotation=None)
6
7 print(" By FactorAnalyzer approach we can see each Factor comprises which stocks\n", '-'*80)
8 print(" Factor1 has AAL,AAPL,ADBE,ADI,ADPADSK... similiar, those having positive values\n", '-'*80)
9 print(fa.loadings)
```

**Problem 2 : How many Unique patterns that exist in the historical stock data set, based on**

**fluctuations in price.**

**Solution 2: Here we do PCA to reduce the dimensions and then KMean Clustering to get Unique patterns in Data**

```
In [ ]: 1 from sklearn.decomposition import PCA
2
3 for ncomponents in range(2,19):
4     pca = PCA(n_components=ncomponents) # two components
5     pca.fit(data) # run PCA, putting in raw version for fun
6     variance = pca.explained_variance_ratio_ #calculate variance ratios
7     var=np.cumsum(np.round(pca.explained_variance_ratio_, decimals=3)*100)
8     print(var) #cumulative sum of variance explained with [n] features
9
10 print(ncomponents)
11 print('-'*80)
12 print("\n The component at which we get maximum variance about 99percent is : ", ncomponents)
```

```
In [ ]: 1 # now for the reduced number of componenets ,apply the clustering
2 pca= PCA(n_components=ncomponents)
3
4 pca.fit(data)
5
6 df1 = pca.transform(data)
7
8
9 print('-'*80,"\n The Stocks Data with the PCA done and reduced Dimensionality\n",'-'*80)
10 print(df1.shape)
11 df2 = pd.DataFrame(df1)
12 print(df2.head())
13
```

```
In [ ]: 1 # using KMeans clustering to know the number of clusters /Unique patterns
2 # Assumption : The Sample size of 10000 taken of data represents the complete Data
3
4 from sklearn.metrics import silhouette_score
5
6 from sklearn.cluster import KMeans
7
8 # Finding the Appropriate number of clusters , when the error is least
9 cluster_range = range( 2, 30 )
10 cluster_errors = []
11
12 X = df2.sample(10000)
13 for num_cluster in cluster_range:
14     cluster = KMeans(n_clusters=num_cluster, random_state=100)
15     cluster.fit(X)
16     error = round(cluster.inertia_,2)
17     cluster_errors.append(error)
18
19 clusters_df = pd.DataFrame( { "num_clusters":cluster_range, "cluster_errors": cluster_errors } )
20
```

```
In [ ]: 1 print(" Following is Elbow method for getting number of clusters with least errors \n", '-'*80)
2 plt.figure(figsize=(12,6))
3 plt.plot( clusters_df.num_clusters, clusters_df.cluster_errors, marker = "o" )
4 plt.ylabel('Custer errors')
5 plt.legend(loc='lower right')
6 plt.title('Elbow Method for number of Clusters Vs Error')
7
8 plt.show()
```

```
In [ ]: 1 # From the Graph of errors we observe that 6 clusters uniquely define the Data
2
3 print('-'*80)
4 print( ''' Based on the Fluctuation in the price, from the error graph we observe\n
5         there are about 6 unique patterns in the Data  ''')
6 print('-'*80)
7 X = df2.sample(10000)
8 print(X.head(2))
9
10 cluster = KMeans(n_clusters=6, random_state=100)
11 cluster.fit_predict(X)
12
13 plt.scatter(X[0],X[1] ,label='True Position')
```

**Problem 3: Identify which all stocks are moving together and which all stocks are different from each other.**

```
In [ ]: 1 # Assumption : 1. Checking for first 10 stocks from the total stocks ,
2 #           same logic to be applied to all the 500 stocks
3 #           2. the threshold for checking if the stocks move together is considered 0.7
4
5
6 # getting the names of all the stocks from the Data Columns
7 import statsmodels.formula.api as sm
8
9 lst=[]
10 diff1st = [[]]
11
12 # List of Stock Names
13 stock = []
14 stockNames = data.columns
15 for col in stockNames:
16     stock.append(str.split(col, '.')[1])
17
18 # for ols apply ols in loop to find the slope and intercept
19
20 #outer loop for the Stock being compared
21 for i in range(15):
22     smstcklst = []
23     diffstcklst = []
24     smstcklst.append(stock[i])
25     diffstcklst.append(stock[i])
26     for j in range(15):
27         if i == j:
28             continue
29         # print(i,j)
30
31         X = data[data.columns[i]]
32         Y = data[data.columns[j]]
33         lm = sm.ols(formula='X ~ Y', data=data)
34         lm = lm.fit()
35         #print(lm.params.Y)
36         if (lm.params.Y > 0.7):
37             smstcklst.append(stock[j])
38         else:
39             diffstcklst.append(stock[j])
40     if(len(smstcklst) != 1):
41         lst.append(smstcklst)
```



```
42     if(len(diffstcklst) != 1):
43         diff1st.append(diffstcklst)
44
45 print(" The Stocks which move together are \n",'-'*80)
46 for i in range(1,len(lst)):
47     print(lst[i])
48 print(" The stocks which are differnt from each other are \n",'-'*80)
49 for i in range(1,len(diff1st)):
50     print(diff1st[i])
51
```