

0.) Import and Clean data

```
In [ ]: import pandas as pd  
# from google.colab import drive  
import matplotlib.pyplot as plt  
import numpy as np  
from sklearn.preprocessing import StandardScaler  
from sklearn.cluster import KMeans
```

```
In [ ]: #drive.mount('/content/gdrive/', force_remount = True)  
df = pd.read_csv("Country-data.csv", sep = ",")
```

```
In [ ]: df.head()
```

```
Out[ ]:
```

	country	child_mort	exports	health	imports	income	inflation	life_expec	total_fer	gdpp
0	Afghanistan	90.2	10.0	7.58	44.9	1610	9.44	56.2	5.82	553
1	Albania	16.6	28.0	6.55	48.6	9930	4.49	76.3	1.65	4090
2	Algeria	27.3	38.4	4.17	31.4	12900	16.10	76.5	2.89	4460
3	Angola	119.0	62.3	2.85	42.9	5900	22.40	60.1	6.16	3530
4	Antigua and Barbuda	10.3	45.5	6.03	58.9	19100	1.44	76.8	2.13	12200

```
In [ ]: y = df['country']  
X = df.drop('country', axis = 1)  
X
```

```
Out[ ]:
```

	child_mort	exports	health	imports	income	inflation	life_expec	total_fer	gdpp
0	90.2	10.0	7.58	44.9	1610	9.44	56.2	5.82	553
1	16.6	28.0	6.55	48.6	9930	4.49	76.3	1.65	4090
2	27.3	38.4	4.17	31.4	12900	16.10	76.5	2.89	4460
3	119.0	62.3	2.85	42.9	5900	22.40	60.1	6.16	3530
4	10.3	45.5	6.03	58.9	19100	1.44	76.8	2.13	12200
...
162	29.2	46.6	5.25	52.7	2950	2.62	63.0	3.50	2970
163	17.1	28.5	4.91	17.6	16500	45.90	75.4	2.47	13500
164	23.3	72.0	6.84	80.2	4490	12.10	73.1	1.95	1310
165	56.3	30.0	5.18	34.4	4480	23.60	67.5	4.67	1310
166	83.1	37.0	5.89	30.9	3280	14.00	52.0	5.40	1460

167 rows × 9 columns

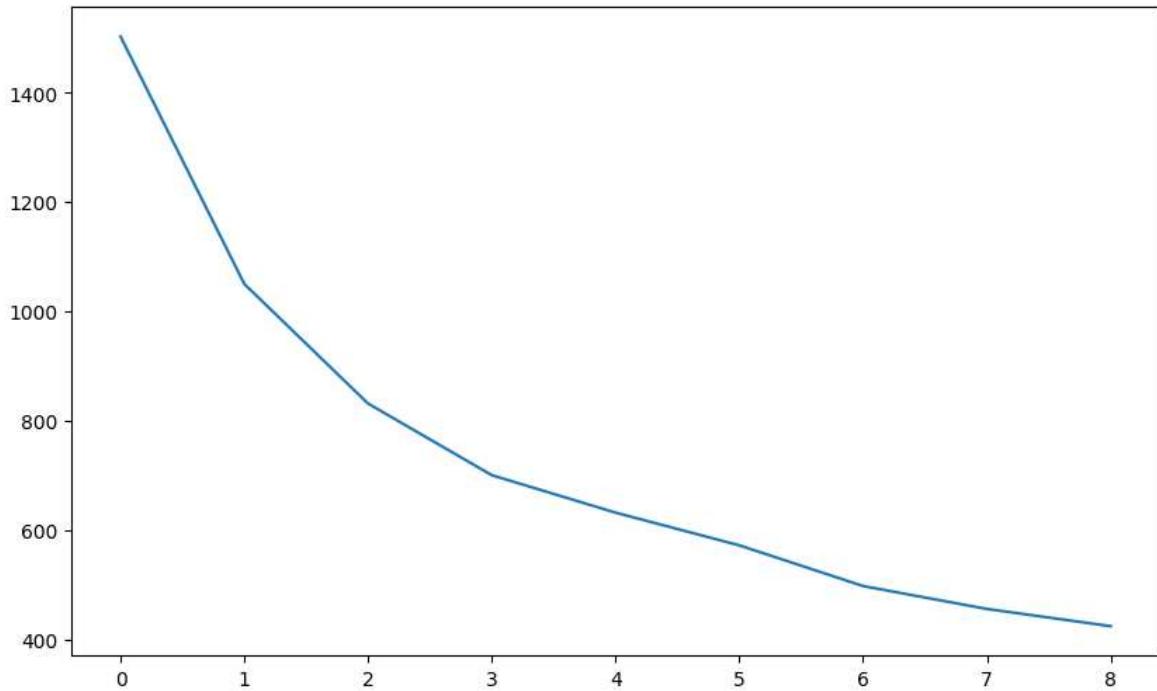
```
In [ ]: names = df[['country']].copy()
```

1.) Fit a kmeans Model with any Number of Clusters

```
In [ ]: scaler = StandardScaler().fit(X)  
X_scaled = scaler.transform(X)
```

```
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:767: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if not hasattr(array, "sparse") and array.dtypes.apply(is_sparse).any():  
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:605: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if is_sparse(pd_dtype):  
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:614: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if is_sparse(pd_dtype) or not is_extension_array_dtype(pd_dtype):  
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:767: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if not hasattr(array, "sparse") and array.dtypes.apply(is_sparse).any():  
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:605: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if is_sparse(pd_dtype):  
c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\validation.py:614: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isinstance(dtype, pd.SparseDtype)` instead.  
if is_sparse(pd_dtype) or not is_extension_array_dtype(pd_dtype):
```

```
In [ ]: kmeans = KMeans(n_clusters = 4)  
kmeans.fit(X_scaled)
```

2.) Pick two features to visualize across

```
In [ ]: X.columns
Out[ ]: Index(['child_mort', 'exports', 'health', 'imports', 'income', 'inflation',
       'life_expec', 'total_fer', 'gdpp'],
       dtype='object')

In [ ]: import matplotlib.pyplot as plt

x1_index = 0
x2_index = 2

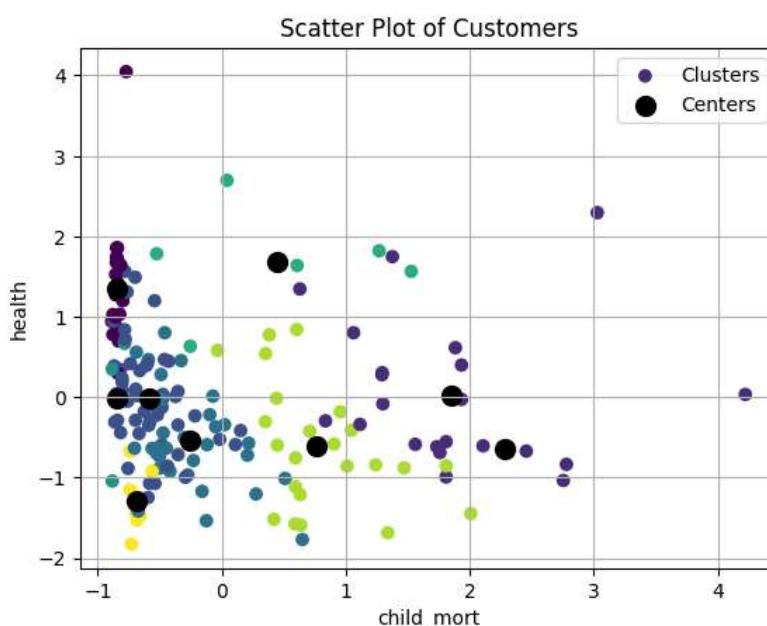
scatter = plt.scatter(X_scaled[:, x1_index], X_scaled[:, x2_index], c=kmeans.labels_, cmap='viridis', label='Clusters')

centers = plt.scatter(kmeans.cluster_centers_[:, x1_index], kmeans.cluster_centers_[:, x2_index], marker='o', color='black', s=100)

plt.xlabel(X.columns[x1_index])
plt.ylabel(X.columns[x2_index])
plt.title('Scatter Plot of Customers')

# Generate legend
plt.legend()

plt.grid()
plt.show()
```



In []:

In []:

3.) Check a range of k-clusters and visualize to find the elbow. Test 30 different random starting places for the centroid means

In []:

```
WCSSs = []
ks = range(1, 15)
for k in ks:
    kmeans = KMeans(n_clusters = k, n_init = 30, init = 'random')
    kmeans.fit(X_scaled)
    WCSSs.append(kmeans.inertia_)
```

```
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
```

In []:

```
# bonus optional do in 1 line of code  
WCSSs = [KMeans(n_clusters = k, n_init = 30, init = 'random').fit(X_scaled).inertia_ for k in range(1, 15)]
```

In [1]: WCSSs # wcsss is lower

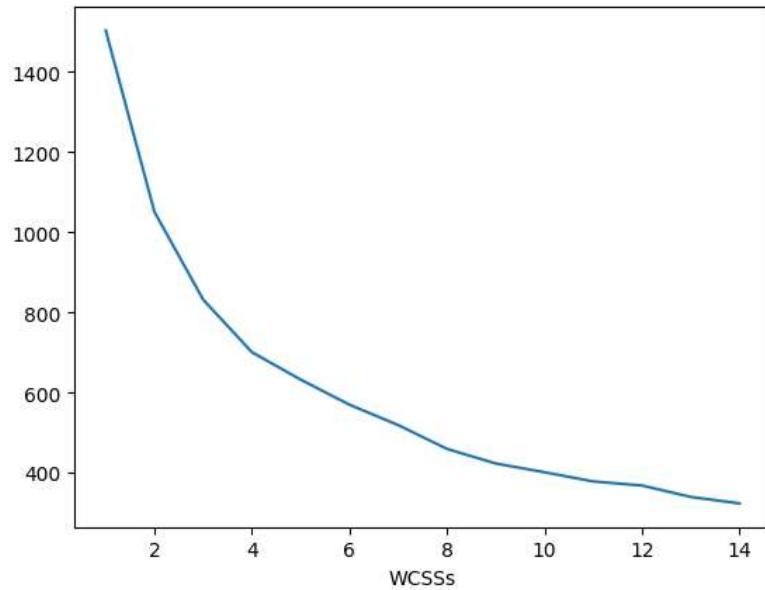
```
Out[ ]: [1502.999999999998,  
        1050.2145582853304,  
        831.4244352086874,  
        700.3229986404376,  
        632.0138132211392,  
        569.5042607233327,  
        518.645818897047,  
        459.2910691015021,  
        422.9155708106426,  
        401.051460603542,  
        378.1104865062954,  
        368.00693560546233,  
        339.46742018928677,  
        323.4035237305541]
```

In []:

4.) Use the above work and economic critical thinking to choose a number of clusters. Explain why you chose the number of clusters and fit a model accordingly.

```
In [ ]: plt.plot(ks, WCSSs)
plt.xlabel('WCSSs')
```

```
Out[1]: Text(0.5, 0, 'WCSSs')
```



6.) Do the same for a silhouette plot

```
In [ ]: from sklearn.metrics import silhouette_score
```

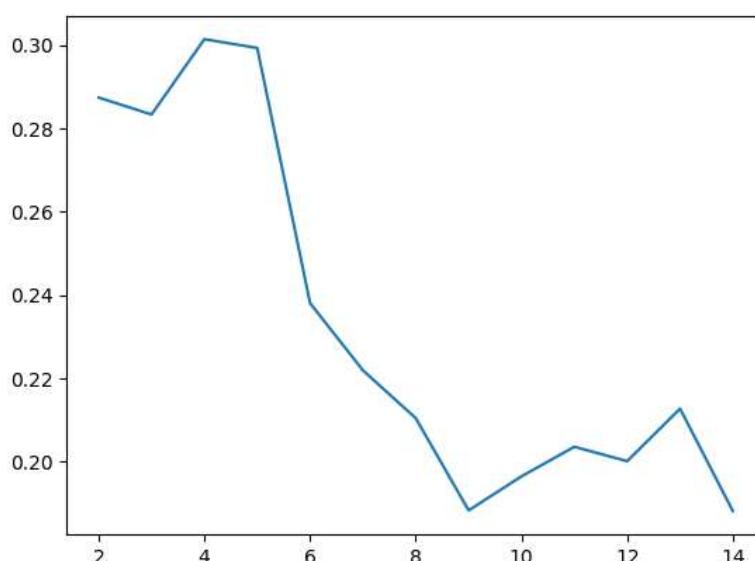


```
In [ ]: CSSs = []
ks = range(2, 15)
for k in ks:
    kmeans = KMeans(n_clusters = k, n_init = 30, init = 'random')
    kmeans.fit(X_scaled)
    labs = kmeans.labels_
    CSSs.append(silhouette_score(X_scaled, labs))
```

```
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
... warnings.warn(
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
```

```
In [ ]: plt.plot(ks, CSSs)
```

```
[out]: []
```



7.) Create a list of the countries that are in each cluster. Write interesting things you notice.

```
In [ ]: kmeans = KMeans(n_clusters = 2, n_init = 30, init = 'random').fit(X_scaled)
```

```
c:\ProgramData\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.  
warnings.warn(
```

```
In [ ]: preds = pd.DataFrame(kmeans.predict(X_scaled))
```

```
In [ ]: output = pd.concat([preds, df], axis = 1)
```

```
In [ ]: output
```

```
Out[ ]:   0      country child_mort exports health imports income inflation life_expec total_fer gdpp
  0  0    Afghanistan     90.2    10.0    7.58    44.9    1610     9.44     56.2     5.82     553
  1  1      Albania     16.6    28.0    6.55    48.6    9930     4.49     76.3     1.65    4090
  2  1      Algeria     27.3    38.4    4.17    31.4   12900    16.10     76.5     2.89    4460
  3  0      Angola    119.0    62.3    2.85    42.9    5900    22.40     60.1     6.16    3530
  4  1  Antigua and Barbuda     10.3    45.5    6.03    58.9   19100     1.44     76.8     2.13   12200
  ... ...
  162 0      Vanuatu     29.2    46.6    5.25    52.7    2950     2.62     63.0     3.50    2970
  163 1      Venezuela     17.1    28.5    4.91    17.6   16500    45.90     75.4     2.47   13500
  164 1      Vietnam     23.3    72.0    6.84    80.2    4490    12.10     73.1     1.95    1310
  165 0      Yemen      56.3    30.0    5.18    34.4    4480    23.60     67.5     4.67    1310
  166 0      Zambia     83.1    37.0    5.89    30.9    3280    14.00     52.0     5.40    1460
```

167 rows × 11 columns

```
In [ ]: output
print('cluster 1:')
list(output.loc[output[0]==0, 'country'])
```

cluster 1:

```
Out[ ]: ['Albania',
 'Algeria',
 'Antigua and Barbuda',
 'Argentina',
 'Armenia',
 'Australia',
 'Austria',
 'Azerbaijan',
 'Bahamas',
 'Bahrain',
 'Barbados',
 'Belarus',
 'Belgium',
 'Belize',
 'Bhutan',
 'Bosnia and Herzegovina',
 'Brazil',
 'Brunei',
 'Bulgaria',
 'Canada',
 'Cape Verde',
 'Chile',
 'China',
 'Colombia',
 'Costa Rica',
 'Croatia',
 'Cyprus',
 'Czech Republic',
 'Denmark',
 'Dominican Republic',
 'Ecuador',
 'El Salvador',
 'Estonia',
 'Fiji',
 'Finland',
 'France',
 'Georgia',
 'Germany',
 'Greece',
 'Grenada',
 'Hungary',
 'Iceland',
 'Iran',
 'Ireland',
 'Israel',
 'Italy',
 'Jamaica',
 'Japan',
 'Jordan',
 'Kazakhstan',
 'Kuwait',
 'Latvia',
 'Lebanon',
 'Libya',
 'Lithuania',
 'Luxembourg',
 'Macedonia, FYR',
 'Malaysia',
 'Maldives',
 'Malta',
 'Mauritius',
 'Moldova',
 'Montenegro',
 'Morocco',
 'Netherlands',
 'New Zealand',
 'Norway',
 'Oman',
 'Panama',
 'Paraguay',
 'Peru',
 'Poland',
 'Portugal',
 'Qatar',
 'Romania',
 'Russia',
 'Saudi Arabia',
 'Serbia',
 'Seychelles',
 'Singapore',
 'Slovak Republic',
 'Slovenia',
 'South Korea',
 'Spain',
 'Sri Lanka',
 'St. Vincent and the Grenadines',
 'Suriname',
 'Sweden',
 'Switzerland',
 'Thailand',
 'Tunisia',
```

```
'Turkey',
'Ukraine',
'United Arab Emirates',
'United Kingdom',
'United States',
'Uruguay',
'Venezuela',
'Vietnam']
```

```
In [ ]: output
print('cluster 2:')
list(output.loc[output[0]==1, 'country'])
```

```
cluster 2:
```

```
Out[ ]: ['Albania',
 'Algeria',
 'Antigua and Barbuda',
 'Argentina',
 'Armenia',
 'Australia',
 'Austria',
 'Azerbaijan',
 'Bahamas',
 'Bahrain',
 'Barbados',
 'Belarus',
 'Belgium',
 'Belize',
 'Bhutan',
 'Bosnia and Herzegovina',
 'Brazil',
 'Brunei',
 'Bulgaria',
 'Canada',
 'Cape Verde',
 'Chile',
 'China',
 'Colombia',
 'Costa Rica',
 'Croatia',
 'Cyprus',
 'Czech Republic',
 'Denmark',
 'Dominican Republic',
 'Ecuador',
 'El Salvador',
 'Estonia',
 'Fiji',
 'Finland',
 'France',
 'Georgia',
 'Germany',
 'Greece',
 'Grenada',
 'Hungary',
 'Iceland',
 'Iran',
 'Ireland',
 'Israel',
 'Italy',
 'Jamaica',
 'Japan',
 'Jordan',
 'Kazakhstan',
 'Kuwait',
 'Latvia',
 'Lebanon',
 'Libya',
 'Lithuania',
 'Luxembourg',
 'Macedonia, FYR',
 'Malaysia',
 'Maldives',
 'Malta',
 'Mauritius',
 'Moldova',
 'Montenegro',
 'Morocco',
 'Netherlands',
 'New Zealand',
 'Norway',
 'Oman',
 'Panama',
 'Paraguay',
 'Peru',
 'Poland',
 'Portugal',
 'Qatar',
 'Romania',
 'Russia',
 'Saudi Arabia',
 'Serbia',
 'Seychelles',
 'Singapore',
 'Slovak Republic',
 'Slovenia',
 'South Korea',
 'Spain',
 'Sri Lanka',
 'St. Vincent and the Grenadines',
 'Suriname',
 'Sweden',
 'Switzerland',
 'Thailand',
 'Tunisia',
```

```
'Turkey',
'Ukraine',
'United Arab Emirates',
'United Kingdom',
'United States',
'Uruguay',
'Venezuela',
'Vietnam']
```

```
In [ ]: ##### Write an observation
```

8.) Create a table of Descriptive Statistics. Rows being the Cluster number and columns being all the features. Values being the mean of the centroid. Use the nonscaled X values for interpretation

```
In [ ]: q8df = pd.concat([preds, X], axis = 1)
```

```
In [ ]: q8df
```

```
Out[ ]:   0  child_mort  exports  health  imports  income  inflation  life_expec  total_fer  gdpp
0  0      90.2     10.0    7.58     44.9    1610     9.44      56.2      5.82     553
1  1      16.6     28.0    6.55     48.6    9930     4.49      76.3     1.65    4090
2  1      27.3     38.4    4.17     31.4   12900     16.10      76.5     2.89    4460
3  0     119.0     62.3    2.85     42.9    5900     22.40      60.1     6.16    3530
4  1      10.3     45.5    6.03     58.9   19100     1.44      76.8     2.13   12200
... ...
162 0      29.2     46.6    5.25     52.7    2950     2.62      63.0     3.50    2970
163 1      17.1     28.5    4.91     17.6   16500     45.90      75.4     2.47   13500
164 1      23.3     72.0    6.84     80.2   4490     12.10      73.1     1.95   1310
165 0      56.3     30.0    5.18     34.4   4480     23.60      67.5     4.67   1310
166 0      83.1     37.0    5.89     30.9   3280     14.00      52.0     5.40   1460
```

167 rows × 10 columns

```
In [ ]: q8df.groupby(0).mean()
```

```
Out[ ]:   child_mort  exports  health  imports  income  inflation  life_expec  total_fer  gdpp
0
0  76.280882  30.198515  6.090147  43.642146  4227.397059  11.098750  61.910294  4.413824  1981.235294
1  12.161616  48.603030  7.314040  49.121212  26017.171717  5.503545  76.493939  1.941111  20507.979798
```

```
In [ ]: q8df.groupby(0).std()
```

```
Out[ ]:   child_mort  exports  health  imports  income  inflation  life_expec  total_fer  gdpp
0
0  38.076068  18.201742  2.645319  19.323451  4890.581414  13.682630  6.897418  1.285590  2528.509189
1  8.523122  30.116032  2.716652  26.928785  20441.749847  6.957187  3.735757  0.486744  20578.727127
```

9.) Write an observation about the descriptive statistics.

- The child mortality of group 0 is higher than that of group 1, on an average level
- The health spending in group 0 is lower, correspondingly
- Group 1 has higher means in economic activities such as exports, income, and GDP. They are more economically developed
- Group 1's life expectancy is higher on average

```
In [ ]:
```