

## 0.) Import and Clean data

```
In [12]: 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 [13]: #drive.mount('/content/gdrive/', force_remount = True)
df = pd.read_csv("Country-data.csv", sep = ",")
```

```
In [14]: df.head()
```

```
Out[14]:
```

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

```
Out[15]:
```

	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

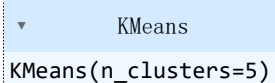
```
In [16]: names = df[['country']].copy()
```

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

```
In [17]: scaler = StandardScaler().fit(X)
X_scaled = scaler.transform(X)
```

```
In [18]: kmeans = KMeans(n_clusters = 5)
kmeans.fit(X_scaled)
```

C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\\_kmeans.py:870: FutureWarning: The default value of `n\_init` will change from 10 to 'auto' in 1.4. Set the value of `n\_init` explicitly to suppress the warning  
warnings.warn(  
C:\Users\Dell\AppData\Roaming\Python\Python311\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(

```
Out[18]: 
KMeans(n_clusters=5)
```

```
In [19]: WCSSs = []      # with in cluster sum of squares
Ks = range(1,10)
for k in Ks:
    kmeans = KMeans(n_clusters = k,n_init=30)
    kmeans.fit(X_scaled)
    WCSSs.append(kmeans.inertia_)
```

```

C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\Dell\AppData\Roaming\Python\Python311\site-packages\sklearn\cluster\_kmeans.py:1382: UserWarn
ing: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than availa
ble threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(

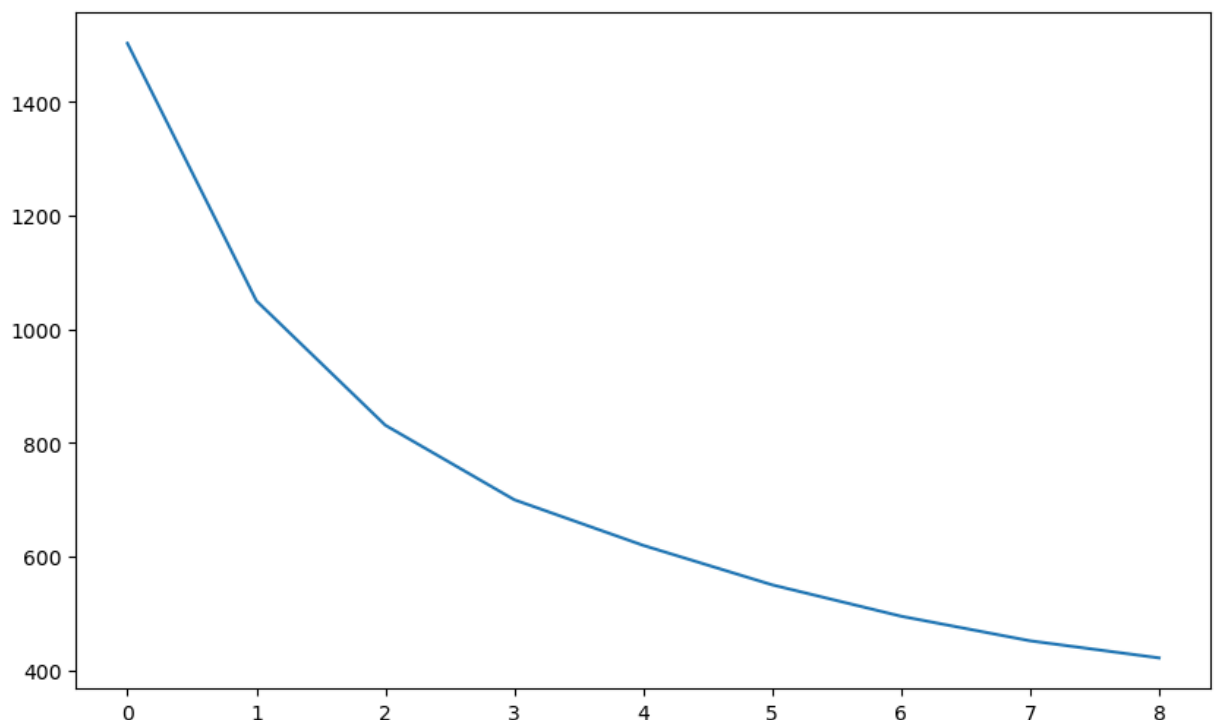
```

In [20]:

```

plt.figure(figsize = (10,6))
plt.plot(WCSSs)
plt.show()

```



## 2.) Pick two features to visualize across

In [21]: `X.columns`

Out[21]: `Index(['child_mort', 'exports', 'health', 'imports', 'income', 'inflation',  
'life_expec', 'total_fer', 'gdpp'],  
 dtype='object')`

```
In [22]: import matplotlib.pyplot as plt

x1_index = 0
x2_index = 2

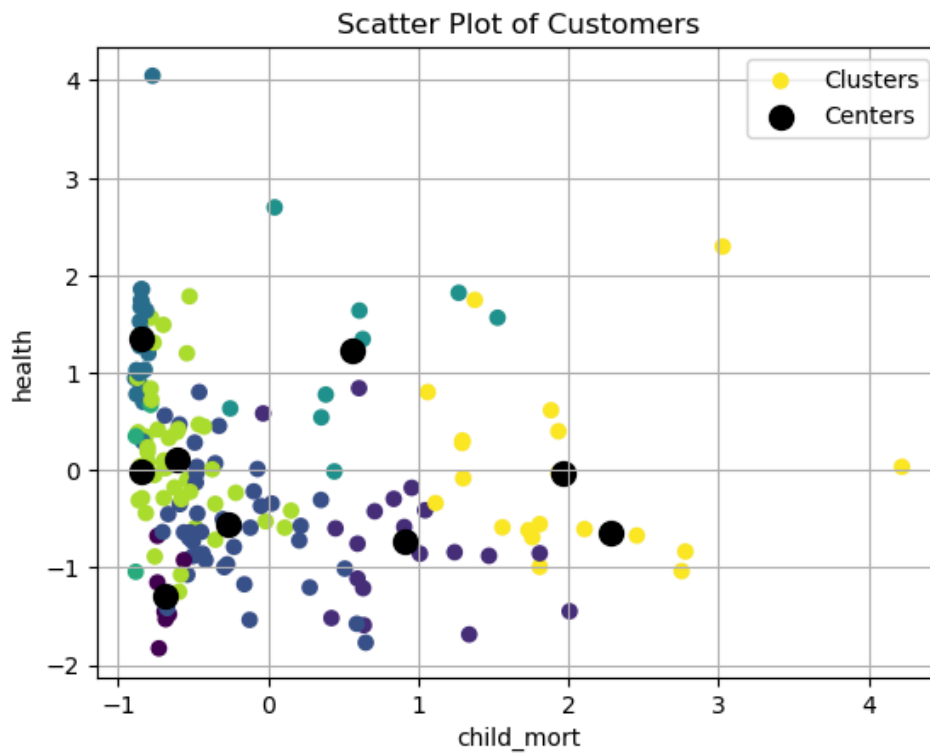
scatter = plt.scatter(X_scaled[:, x1_index], X_scaled[:, x2_index], c=kmeans.labels_,
                      cmap='viridis', label='Clusters')

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

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()
```



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

```
In [23]: 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_)
```



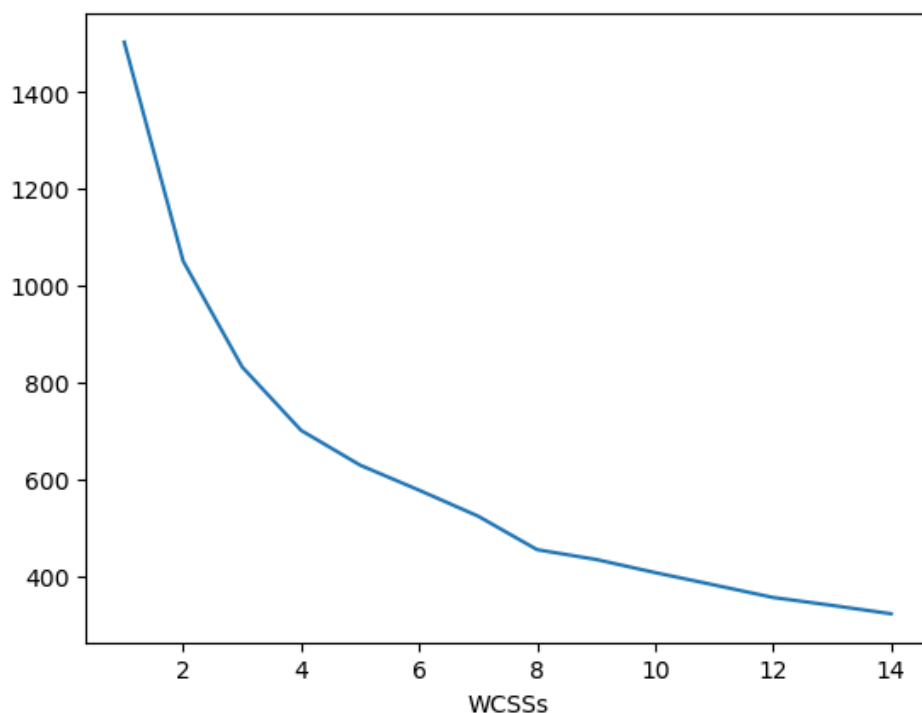


```
Out[25]: [1503.0000000000002,  
1050.2145582853304,  
831.4244352086873,  
700.3229986404375,  
628.5382539105315,  
576.7432043038536,  
523.1157817835406,  
453.84219458741256,  
433.5868887904837,  
406.4888459654401,  
380.9516144118544,  
355.06539455450934,  
338.6653859589953,  
321.37523632099794]
```

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

```
Out[26]: Text(0.5, 0, 'WCSSs')
```



**6.) Do the same for a silhouette plot**

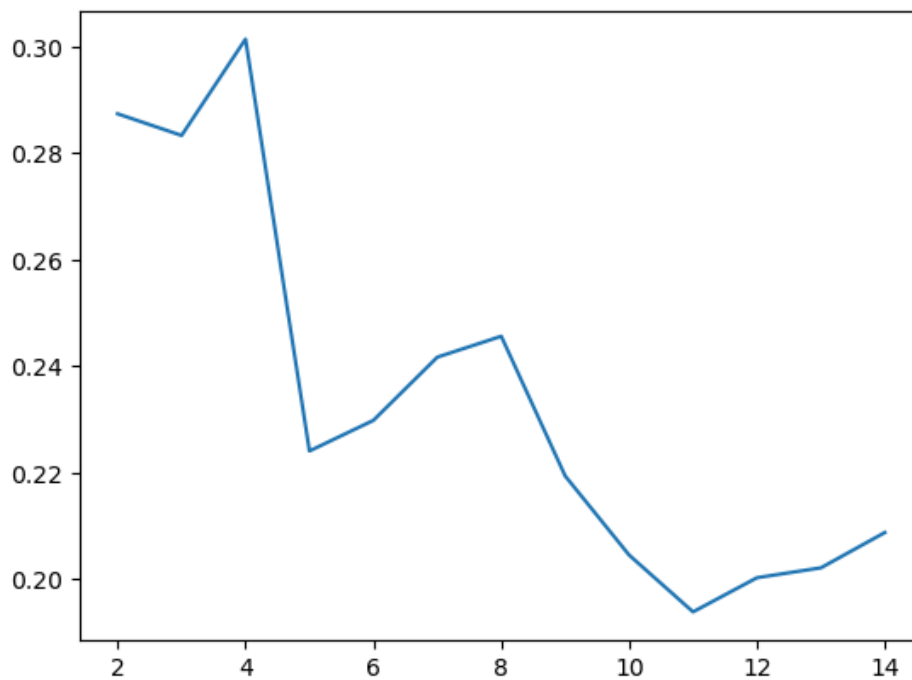
```
In [27]: from sklearn.metrics import silhouette_score
```

```
In [28]: silhouette_score?
```

```
In [29]: CSSs = []  
ks = range(2, 15)
```







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

```
In [31]: kmeans = KMeans(n_clusters = 2, n_init = 30, init = 'random').fit(X_scaled)
```

C:\Users\Dell\AppData\Roaming\Python\Python311\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 [32]: preds = pd.DataFrame(kmeans.predict(X_scaled))
```

```
In [33]: output = pd.concat([preds, df], axis = 1)
```

```
In [34]: output
```

```
Out[34]:
```

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

167 rows × 11 columns

In [35]:

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

cluster 1:

Out[35]:

```
['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 [58]:

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

```
cluster 2:
```

```
Out[58]: ['Afghanistan',
          'Angola',
          'Bangladesh',
          'Benin',
          'Bolivia',
          'Botswana',
          'Burkina Faso',
          'Burundi',
          'Cambodia',
          'Cameroon',
          'Central African Republic',
          'Chad',
          'Comoros',
          'Congo, Dem. Rep.',
          'Congo, Rep.',
          "Cote d'Ivoire",
          'Egypt',
          'Equatorial Guinea',
          'Eritrea',
          'Gabon',
          'Gambia',
          'Ghana',
          'Guatemala',
          'Guinea',
          'Guinea-Bissau',
          'Guyana',
          'Haiti',
          'India',
          'Indonesia',
          'Iraq',
          'Kenya',
          'Kiribati',
          'Kyrgyz Republic',
          'Lao',
          'Lesotho',
          'Liberia',
          'Madagascar',
          'Malawi',
          'Mali',
          'Mauritania',
          'Micronesia, Fed. Sts.',
          'Mongolia',
          'Mozambique',
          'Myanmar',
          'Namibia',
          'Nepal',
          'Niger',
          'Nigeria',
          'Pakistan',
          'Philippines',
          'Rwanda',
          'Samoa',
          'Senegal',
          'Sierra Leone',
          'Solomon Islands',
          'South Africa',
          'Sudan',
          'Tajikistan',
          'Tanzania',
          'Timor-Leste',
          'Togo',
          'Tonga',
          'Turkmenistan',
          'Uganda',
          'Uzbekistan',
          'Vanuatu',
```

```
'Yemen',  
'Zambia']
```

```
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 interprotation

```
In [36]: q8df = pd.concat([preds, X], axis = 1)
```

```
In [37]: q8df
```

```
Out[37]:
```

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

167 rows × 10 columns

```
In [38]: q8df.groupby(0).mean()
```

```
Out[38]:
```

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

9.) Write an observation about the descriptive statistics.

Group 0, on average, has a higher child mortality rate and total fertility rate than Group 1. Group 1 has higher means for exports, income, and GDP per capita, which could suggest it consists of more economically developed countries or entities compared to Group 0. The health spending is higher on average in Group 1 compared to Group 0. Life expectancy is higher on average in Group 1 than in Group 0, which could be correlated with the group's higher income and health spending. The inflation mean is higher in Group 0 than in Group 1, which might reflect economic challenges faced by the countries in Group 0.

