


```
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/gdp-per-capita-maddison.csv')

# Define list of EU member states by their ISO country codes
eu_countries = [
    'BEL', 'BGR', 'CZE', 'DNK', 'DEU', 'EST', 'IRL', 'GRC', 'ESP', 'FRA',
    'HRV', 'ITA', 'CYP', 'LVA', 'LTU', 'LUX', 'HUN', 'MLT', 'NLD', 'AUT',
    'POL', 'PRT', 'ROU', 'SVN', 'SVK', 'FIN', 'SWE'
]

# Filter dataset for years 2000 to 2018
df_2000_2018 = df[(df['Year'] >= 2000) & (df['Year'] <= 2018) & (df['Code'].isin(eu_countries))]

# Pivot table to have countries as rows and years as columns
pivot_table = df_2000_2018.pivot(index='Code', columns='Year', values='GDP per capita')

# Sort table by last year to see the progression
pivot_table_sorted = pivot_table.sort_values(by=2018, ascending=False)

# Drop any missing values if present (countries with missing data)
pivot_table_cleaned = pivot_table_sorted.dropna()

# Standardize data (important for k-means)
scaler = StandardScaler()
data_scaled = scaler.fit_transform(pivot_table_cleaned)

# Elbow Method to determine k
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    kmeans.fit(data_scaled)
    wcss.append(kmeans.inertia_)

# Plot Elbow graph
plt.figure(figsize=(8, 4))
plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS') # WCSS stands for Within-Cluster Sum of Square
plt.show()

# Based on Elbow graph, choose number of clusters (k)
k = 4
kmeans = KMeans(n_clusters=k, init='k-means++', max_iter=300, n_init=10, random_state=0)

# Fit KMeans using standardized data
clusters = kmeans.fit_predict(data_scaled)

# Perform PCA to reduce dimensions to 2 for visualization
pca = PCA(n_components=2)
principal_components = pca.fit_transform(data_scaled)

# Create a new DataFrame for PCA results
pca_df = pd.DataFrame(data=principal_components, columns=['PC1', 'PC2'])

# Add cluster labels to PCA DataFrame
pca_df['Cluster'] = clusters

# Add country codes to PCA DataFrame for labeling
pca_df['Country'] = pivot_table_cleaned.index

# Get centroids
centroids = kmeans.cluster_centers_

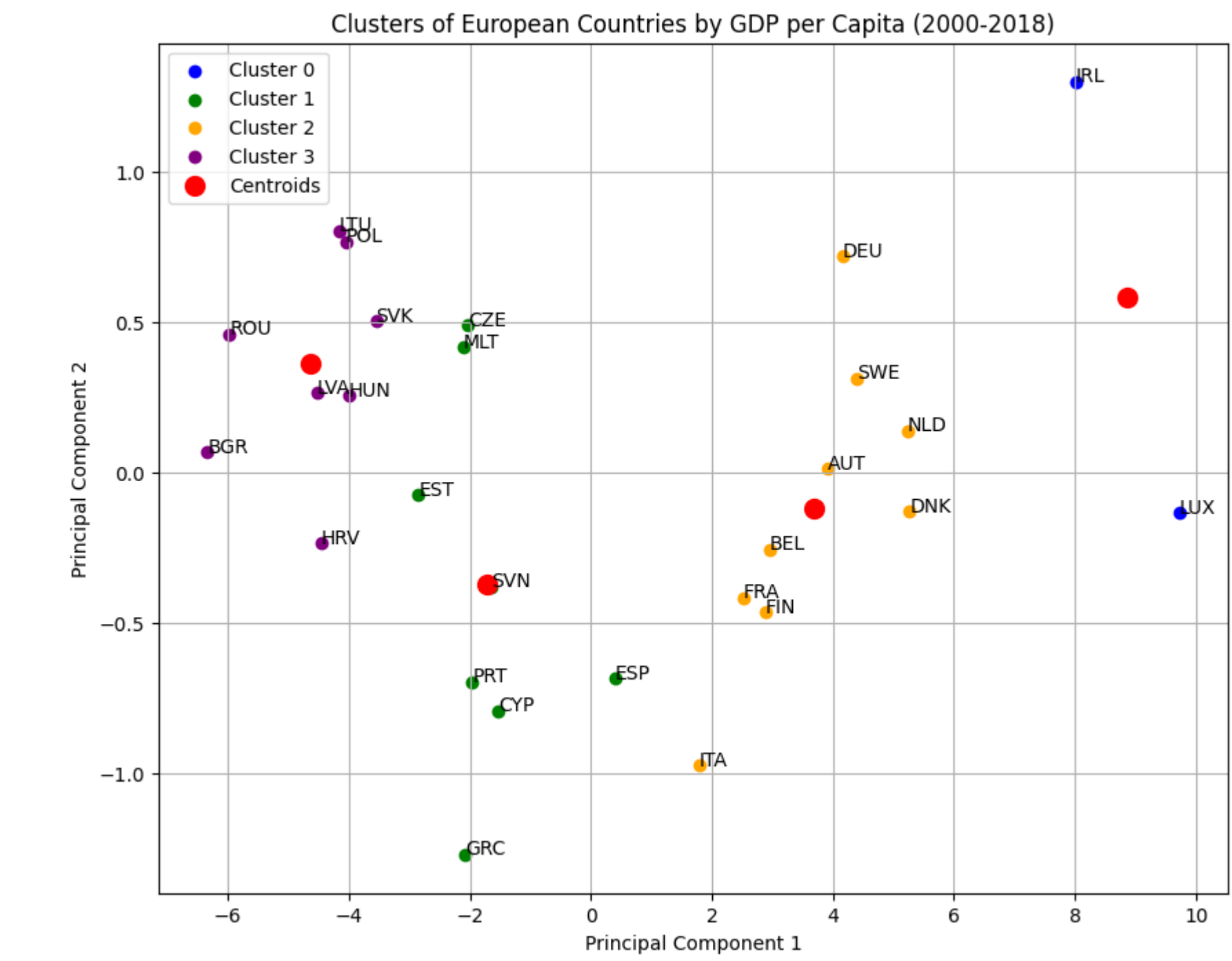
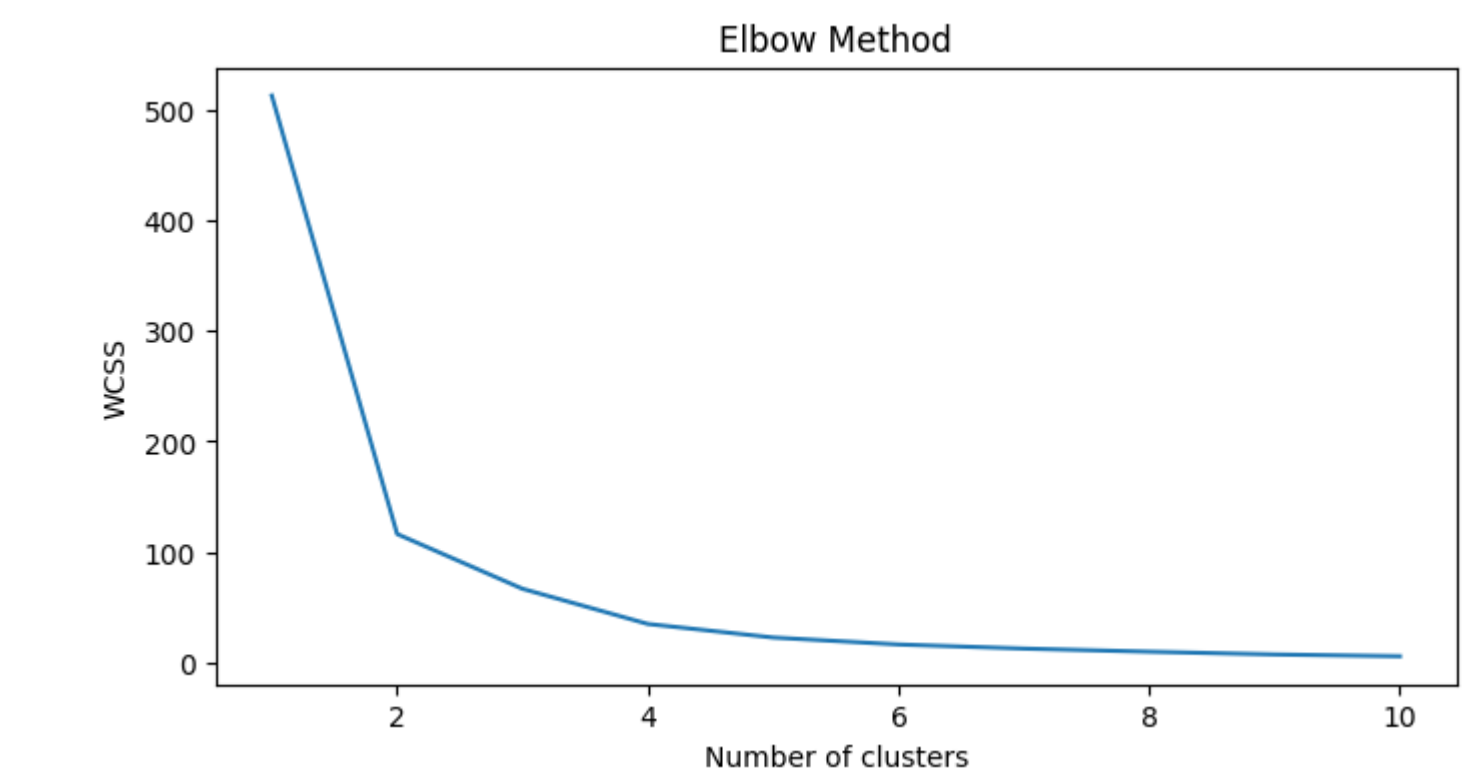
# Transform centroids using PCA model
centroids_pca = pca.transform(centroids)

# Plot clusters
plt.figure(figsize=(10, 8))
colors = ['blue', 'green', 'orange', 'purple']
for i in range(k):
    plt.scatter(pca_df[pca_df['Cluster'] == i]['PC1'], pca_df[pca_df['Cluster'] == i]['PC2'], label=f'Cluster {i}', c=colors[i])

# Plotting centroids
plt.scatter(centroids_pca[:, 0], centroids_pca[:, 1], s=100, c='red', label='Centroids')

plt.title('Clusters of European Countries by GDP per Capita (2000-2018)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.legend()
plt.grid(True)

# Annotate country codes
for i, txt in enumerate(pca_df['Country']):
    plt.annotate(txt, (pca_df['PC1'][i], pca_df['PC2'][i]))
plt.show()
```



```
import pandas as pd
import numpy as np
from sklearn.impute import KNNImputer

df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Gener of plastic packaging waste per capita.csv', skiprows=8)

# only need first 20 columns, ignore rest
df = df.iloc[:, :20] # discard empty columns

# Rename columns assuming first column is 'GEO (Labels)' and rest are years from 2000 to 2018
df.columns = ['Code'] + list(range(2000, 2019))

# Convert all entries in 'GEO (Labels)' to uppercase to match country codes in eu_countries list
df['Code'] = df['Code'].str.upper()

eu_countries = [
    'BEL', 'BGR', 'CZE', 'DNK', 'DEU', 'EST', 'IRL', 'GRC', 'ESP', 'FRA',
    'HRV', 'ITA', 'CYP', 'LVA', 'LTU', 'LUX', 'HUN', 'MLT', 'NLD', 'AUT',
    'POL', 'PRT', 'ROU', 'SVN', 'SVK', 'FIN', 'SWE'
]

# Filter dataframe to include only rows where 'GEO (Labels)' is in eu_countries list
eu_data = df[df['Code'].isin(eu_countries)]

# Replace non-numeric values with NaN and convert all columns to numeric, coercing errors to NaN
for column in eu_data.columns[1:]: # Skipping 'GEO (Labels)' column
    eu_data.loc[:, column] = pd.to_numeric(eu_data[column], errors='coerce')

# Limit data to years 2000 to 2018
eu_data = eu_data[['Code'] + list(range(2000, 2019))]

# Display cleaned EU data for years 2000-2018
print(eu_data)
```

	2010	2011	2012	2013	2014	2015	2016	2017	2018
0	41.14	34.65	36.65	44.40	59.32	60.01	57.94	58.38	54.24
1	43.96	45.92	45.73	50.10	50.02	46.88	46.14	46.39	42.61
2	38.12	39.18	35.98	49.06	50.17	46.47	49.10	49.95	41.90
3	32.90	34.58	35.27	35.63	36.37	37.36	37.62	38.53	39.83
4	34.14	33.79	33.31	34.16	34.59	35.70	36.66	38.86	40.31
5	34.04	34.94	34.46	33.91	34.25	35.05	36.53	37.52	37.93
6	29.82	33.80	32.85	33.85	33.22	34.67	37.46	34.81	42.88
7	31.63	31.48	32.24	34.05	34.16	34.12	34.09	34.36	34.16
8	30.88	31.20	30.53	30.09	31.10	32.06	32.65	34.80	35.09
9	30.01	28.99	27.89	28.00	30.52	31.75	32.84	34.53	35.37
10	21.09	20.93	25.90	27.85	26.21	30.46	31.48	32.24	34.84
11	28.96	28.62	28.85	29.51	29.36	30.13	30.28	30.36	30.40
12	27.32	26.60	27.39	27.85	28.11	29.04	29.54	29.89	30.35
13	29.34	27.39	25.82	26.95	25.72	28.03	31.94	28.41	31.80
14	19.27	20.61	21.86	23.53	23.60	24.63	26.53	27.42	25.94
15	21.16	22.43	22.44	23.18	23.55	23.57	24.03	23.93	24.17
16	20.01	19.95	20.14	20.40	20.79	23.45	22.42	23.59	25.16
17	18.25	19.93	19.98	21.37	22.87	22.55	22.88	25.42	27.08
18	22.10	21.79	21.80	20.41	21.44	21.85	22.45	24.28	23.81
19	21.67	21.74	21.05	21.65	21.38	21.27	22.96	23.66	24.52
20	16.78	17.57	18.18	20.14	19.41	20.92	20.55	20.31	22.63
21	19.62	19.75	19.33	18.06	18.03	19.62	21.99	22.83	24.22

A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
eu_data.loc[:, column] = pd.to_numeric(eu_data[column], errors='coerce')
<ipython-input-4-ac213def497e>:28: DeprecationWarning: In a future version, 'df.iloc[:, i] = newvals' will attempt to set the values inplace instead of always setting a new array. To retain the old behavior, use either 'df[df.columns[i]] = newvals' or, if columns are non-unique, 'df.isetitem(i, newvals)'

```
from sklearn.impute import KNNImputer
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import pandas as pd

# Convert year columns to integers if they are not already
year_columns = list(range(2000, 2019))

# Apply KNN imputer to these columns
imputer = KNNImputer(n_neighbors=5)
eu_data[year_columns] = imputer.fit_transform(eu_data[year_columns])

# Use Elbow Method to find optimal number of clusters
wcss = []
for i in range(1, 11): # Test 1 to 10 clusters or adjust range as needed
    kmeans = KMeans(n_clusters=i, random_state=0)
    kmeans.fit(eu_data[year_columns])
    wcss.append(kmeans.inertia_)

# Plot Elbow graph
plt.figure(figsize=(8, 6))
plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()

# select optimal number of clusters based on plot
optimal_clusters = 4

# Fit KMeans model with optimal number of clusters
kmeans = KMeans(n_clusters=optimal_clusters, random_state=0)
eu_data['cluster'] = kmeans.fit_predict(eu_data[year_columns])

# Perform PCA to reduce data to 2 dimensions for visualization
pca = PCA(n_components=2)
reduced_data = pca.fit_transform(eu_data[year_columns])

# Get cluster assignments and country codes
clusters = eu_data['cluster'].values
country_codes = eu_data['Code'].values # Assuming 'GEO (Labels)' is column with country codes

# Scatter plot of reduced data with cluster assignments
plt.figure(figsize=(12, 10))
scatter = plt.scatter(reduced_data[:, 0], reduced_data[:, 1], c=clusters, cmap='viridis', alpha=0.6)

# Annotate each data point with country code
for i, txt in enumerate(country_codes):
    plt.annotate(txt, (reduced_data[i, 0], reduced_data[i, 1]), fontsize=9)

# Plotting centroids (transformed with PCA)
centroids = pca.transform(kmeans.cluster_centers_)
plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=200, c='red', label='Centroids')

# Adding labels and title
plt.xlabel('PCA Feature 1')
plt.ylabel('PCA Feature 2')
plt.title('2D PCA of EU Countries Clustering')

# Adding legend for clusters
plt.legend(*scatter.legend_elements(), title='Clusters')

plt.grid(True)
plt.tight_layout()

plt.show()

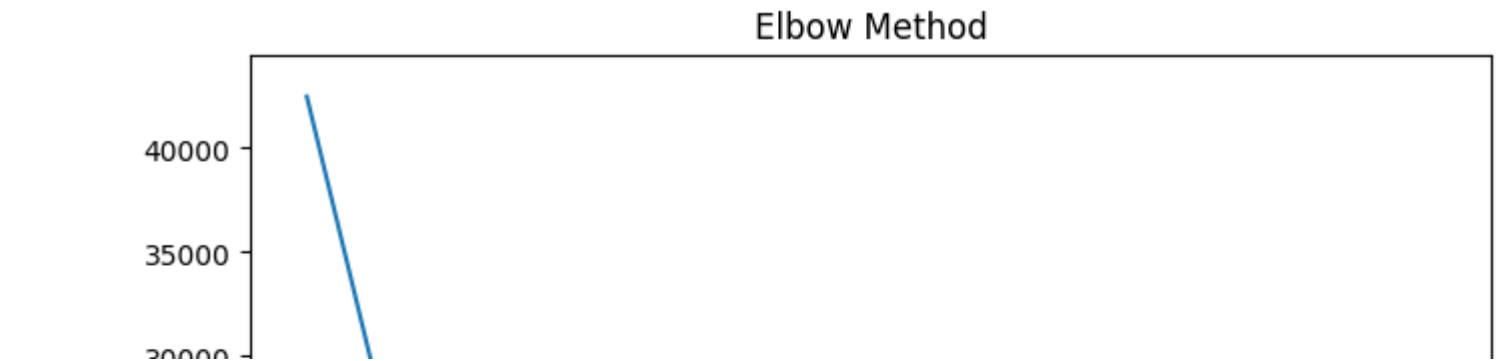
# Print DataFrame with imputed values and cluster assignments
print(eu_data[['Code'] + year_columns + ['cluster']])
```



```

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(

```



```

# Reset index to make sure 'Code' is a column in eu_data
eu_data.reset_index(drop=True, inplace=True)

```

```

# merge pivot_table_sorted and eu_data on 'Code' column
merged_data = pd.merge(pivot_table_sorted, eu_data, on='Code')
# Remove 'cluster' column from dataframe
merged_data.drop('cluster', axis=1, inplace=True)

```

```

#rename columns for clarity
column_mapping = {
    '2000_x': 'GDP_2000', '2001_x': 'GDP_2001', '2002_x': 'GDP_2002', '2003_x': 'GDP_2003',
    '2004_x': 'GDP_2004', '2005_x': 'GDP_2005', '2006_x': 'GDP_2006', '2007_x': 'GDP_2007',
    '2008_x': 'GDP_2008', '2009_x': 'GDP_2009', '2010_x': 'GDP_2010', '2011_x': 'GDP_2011',
    '2012_x': 'GDP_2012', '2013_x': 'GDP_2013', '2014_x': 'GDP_2014', '2015_x': 'GDP_2015',
    '2016_x': 'GDP_2016', '2017_x': 'GDP_2017', '2018_x': 'GDP_2018',
    '2000_y': 'Plastic_Waste_2000', '2001_y': 'Plastic_Waste_2001', '2002_y': 'Plastic_Waste_2002',
    '2003_y': 'Plastic_Waste_2003', '2004_y': 'Plastic_Waste_2004', '2005_y': 'Plastic_Waste_2005',
    '2006_y': 'Plastic_Waste_2006', '2007_y': 'Plastic_Waste_2007', '2008_y': 'Plastic_Waste_2008',
    '2009_y': 'Plastic_Waste_2009', '2010_y': 'Plastic_Waste_2010', '2011_y': 'Plastic_Waste_2011',
    '2012_y': 'Plastic_Waste_2012', '2013_y': 'Plastic_Waste_2013', '2014_y': 'Plastic_Waste_2014',
    '2015_y': 'Plastic_Waste_2015', '2016_y': 'Plastic_Waste_2016', '2017_y': 'Plastic_Waste_2017',
    '2018_y': 'Plastic_Waste_2018'
}

```

```

merged_data.rename(columns=column_mapping, inplace=True)
pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)
print(merged_data)
merged_data.to_csv('content/drive/MyDrive/Colab Notebooks/merged_data.csv', index=False)

```

```

2      27.39      27.85      28.11
3      32.85      33.85      33.22
4      35.27      35.63      36.37
5      22.44      23.18      23.55
6      32.24      34.05      34.16
7      28.85      29.51      29.36
8      21.65      21.65      21.38
9      30.53      30.09      31.10
10     34.46      33.91      34.25
11     25.82      26.95      25.72
12     27.89      28.00      30.52
13     20.14      20.46      20.79
14     21.80      20.41      21.44
15     21.86      23.53      23.60
16     35.98      49.06      50.17
17     19.98      21.37      22.87
18     17.62      18.25      18.56
19     19.33      18.06      18.03
20     33.31      34.16      34.59
21     25.90      27.85      26.21
22     18.18      20.14      19.41
23     16.74      16.55      16.93
24     11.31      11.46      11.59
25     14.86      14.53      16.92

Plastic_Waste_2015 Plastic_Waste_2016 Plastic_Waste_2017 \
0      60.81      57.94      58.38
1      46.88      46.14      46.39
2      29.84      29.54      29.89
3      34.67      37.46      34.81
4      37.36      37.62      38.53
5      23.57      24.03      23.93
6      34.12      34.09      34.36
7      30.13      30.28      30.36
8      21.27      22.36      23.66
9      32.06      32.65      34.80
10     35.05      36.53      37.52
11     28.03      31.94      28.41
12     31.75      32.84      34.53
13     23.45      22.42      23.59
14     21.85      22.45      24.28
15     24.63      26.53      27.42
16     46.47      49.10      49.95
17     22.55      22.88      25.42
18     18.99      19.34      20.47
19     19.62      21.99      22.83
20     35.70      36.66      38.86
21     30.46      31.48      32.24
22     20.92      20.55      20.31
23     16.99      17.32      17.50
24     12.35      13.12      14.67
25     18.12      17.70      18.40

Plastic_Waste_2018
0      54.24
1      42.61
2      30.35
3      42.88

```

```

# Filter DataFrame to get two separate DataFrames for GDP and Plastic Waste
gdp_columns = [col for col in merged_data.columns if 'GDP' in col]
plastic_waste_columns = [col for col in merged_data.columns if 'Plastic_Waste' in col]

# Get DataFrame for GDP and Plastic Waste
gdp_data = merged_data[['Code'] + gdp_columns]
plastic_waste_data = merged_data[['Code'] + plastic_waste_columns]

#calculate summary statistics for GDP for each country
gdp_stats = gdp_data.set_index('Code').stack().groupby('Code').agg(['mean', 'min', 'max'])

#calculate summary statistics for Plastic Waste for each country
plastic_waste_stats = plastic_waste_data.set_index('Code').stack().groupby('Code').agg(['mean', 'min', 'max'])

```

```

> 63.720 ... 38.220 39.180 35.98 49.06 50.17 46.47 49.10 49.95

import matplotlib.pyplot as plt

# Filter DataFrame to get two separate DataFrames for GDP and Plastic Waste
gdp_columns = [col for col in merged_data.columns if 'GDP' in col]
plastic_waste_columns = [col for col in merged_data.columns if 'Plastic_Waste' in col]

# Get the DataFrame for GDP and Plastic Waste
gdp_data = merged_data[['Code'] + gdp_columns]
plastic_waste_data = merged_data[['Code'] + plastic_waste_columns]

#calculate summary statistics for GDP for each country
gdp_stats = gdp_data.set_index('Code').stack().groupby('Code').agg(['mean', 'min', 'max'])

#calculate summary statistics for Plastic Waste for each country
plastic_waste_stats = plastic_waste_data.set_index('Code').stack().groupby('Code').agg(['mean', 'min', 'max'])

# Function to create a bar chart for statistics
def plot_stats(df, title):
    # Create a figure and a set of subplots
    fig, ax = plt.subplots(figsize=(10, 6))

    # Plot mean, min, and max
    df['mean'].plot(kind='bar', ax=ax, color='skyblue', position=0, label='Mean')
    df['min'].plot(kind='bar', ax=ax, color='lightgreen', position=1, label='Min')
    df['max'].plot(kind='bar', ax=ax, color='salmon', position=2, label='Max')

    # Set title and labels
    ax.set_title(title)
    ax.set_ylabel('Value')
    ax.set_xlabel('Country Code')
    ax.legend()

    # Show plot
    plt.xticks(rotation=90) # Rotate x-axis labels to show them better
    plt.tight_layout()
    plt.show()

# Plot GDP stats
plot_stats(gdp_stats, "GDP Statistics by Country")

# Plot Plastic Waste stats
plot_stats(plastic_waste_stats, "Plastic Waste Statistics by Country")

```

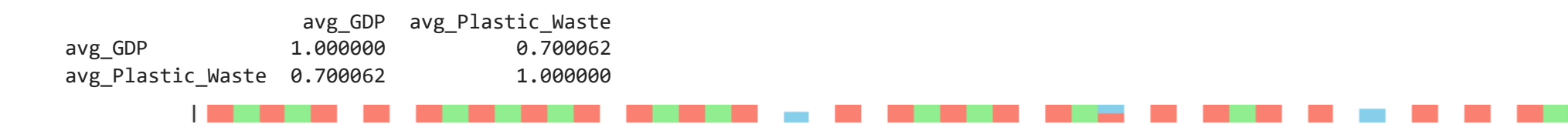
GDP Data by Country

```
import pandas as pd

# Calculate average of GDP and Plastic Waste over years for each country
merged_data['avg_GDP'] = merged_data[gdp_columns].mean(axis=1)
merged_data['avg_Plastic_Waste'] = merged_data[plastic_waste_columns].mean(axis=1)

# Use .corr() method to find Pearson correlation coefficient
correlation_matrix = merged_data[['avg_GDP', 'avg_Plastic_Waste']].corr()

# Show correlation matrix
print(correlation_matrix)
```



```
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt

# Filter DataFrame to only include GDP and Plastic Waste columns
features = merged_data[['avg_GDP', 'avg_Plastic_Waste']]

# Standardize features
scaler = StandardScaler()
features_scaled = scaler.fit_transform(features)

# Choose number of clusters (k) and fit KMeans model
k = 4
kmeans = KMeans(n_clusters=k, random_state=42)
kmeans.fit(features_scaled)

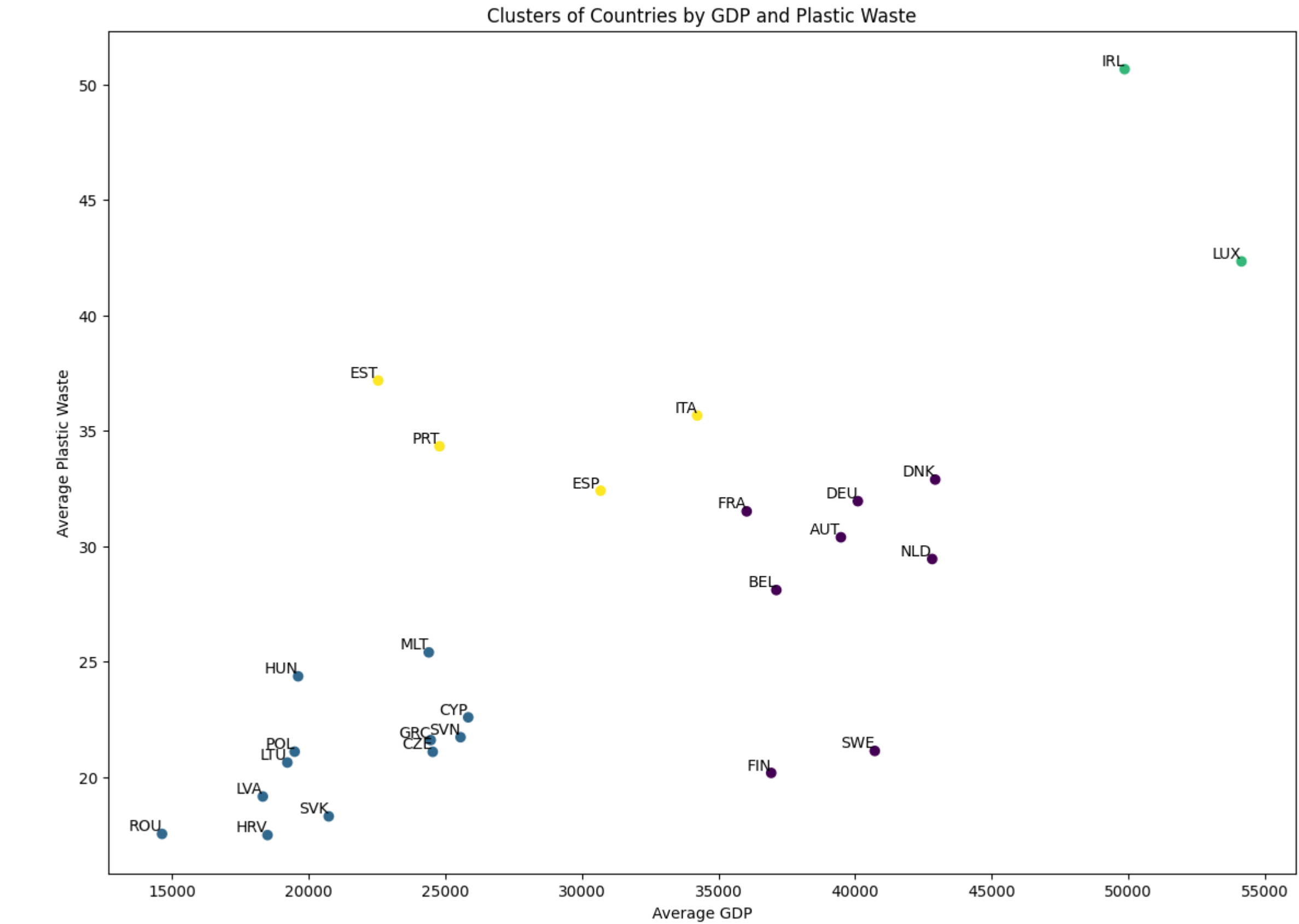
# Add cluster information back to original DataFrame
merged_data['cluster'] = kmeans.labels_

plt.figure(figsize=(14, 10))
plt.scatter(merged_data['avg_GDP'], merged_data['avg_Plastic_Waste'], c=merged_data['cluster'], cmap='viridis')

# Annotate each point in the scatter plot with country code
for i, row in merged_data.iterrows():
    plt.text(row['avg_GDP'], row['avg_Plastic_Waste'], row['Code'], color='black', ha='right', va='bottom')

plt.title('Clusters of Countries by GDP and Plastic Waste')
plt.xlabel('Average GDP')
plt.ylabel('Average Plastic Waste')
plt.show()
```

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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(



```
import pandas as pd

file_path = '/content/drive/MyDrive/Colab Notebooks/globalpop.xls'

global_df = pd.read_excel(file_path, header=3)
eu_df = global_df[global_df['Country Code'].isin(eu_countries)]
years = [str(year) for year in range(2000, 2019)] # Years from 2000 to 2018
columns_to_keep = ['Country Name', 'Country Code', 'Indicator Name', 'Indicator Code'] + years
eu_df_years = eu_df[columns_to_keep]

years = range(2000, 2019)
for year in years:
    eu_df_years.rename(columns={str(year): f'Population_{year}'}, inplace=True)

# Import pandas as pd
file_path = '/content/drive/MyDrive/Colab Notebooks/merged_data.csv'
merged_data = pd.read_csv(file_path)

combined_df = pd.merge(eu_df_years, merged_data, how='inner', left_on='Country Code', right_on='Code')
output_file_path = '/content/drive/MyDrive/Colab Notebooks/combined_data.csv'
combined_df.to_csv(output_file_path, index=False)
print(combined_df.head(20))
```

	Country Name	Country Code	Indicator Name	Indicator Code	\
0	Austria	AUT	Population, total	SP.POP.TOTL	
1	Belgium	BEL	Population, total	SP.POP.TOTL	
2	Cyprus	CYP	Population, total	SP.POP.TOTL	
3	Czechia	CZE	Population, total	SP.POP.TOTL	
4	Germany	DEU	Population, total	SP.POP.TOTL	
5	Denmark	DNK	Population, total	SP.POP.TOTL	
6	Spain	ESP	Population, total	SP.POP.TOTL	
7	Estonia	EST	Population, total	SP.POP.TOTL	
8	Finland	FIN	Population, total	SP.POP.TOTL	
9	France	FRA	Population, total	SP.POP.TOTL	
10	Greece	GRC	Population, total	SP.POP.TOTL	
11	Croatia	HRV	Population, total	SP.POP.TOTL	
12	Hungary	HUN	Population, total	SP.POP.TOTL	
13	Ireland	IRL	Population, total	SP.POP.TOTL	
14	Italy	ITA	Population, total	SP.POP.TOTL	
15	Lithuania	LTU	Population, total	SP.POP.TOTL	
16	Luxembourg	LUX	Population, total	SP.POP.TOTL	
17	Latvia	LVA	Population, total	SP.POP.TOTL	
18	Malta	MLT	Population, total	SP.POP.TOTL	
19	Netherlands	NLD	Population, total	SP.POP.TOTL	

	Population_2000	Population_2001	Population_2002	Population_2003	\
0	8011566.0	8042293.0	8081957.0	8121423.0	
1	10251258.0	10286570.0	10332785.0	10376133.0	
2	946237.0	964830.0	982194.0	1000358.0	
3	10255063.0	10216605.0	10196916.0	10193998.0	
4	82211508.0	82349925.0	82488495.0	82534176.0	
5	5339616.0	5358763.0	5375931.0	5390574.0	
6	40567864.0	40850412.0	41431558.0	42187645.0	
7	1396985.0	1388115.0	1379350.0	1370720.0	
8	5176209.0	5189088.0	5200598.0	5213014.0	
9	60921384.0	61367388.0	61816234.0	62256970.0	
10	10805808.0	10862132.0	10902022.0	10928070.0	
11	4468302.0	4299642.0	4302174.0	4303399.0	
12	10210971.0	10187576.0	10158608.0	10129552.0	
13	3805374.0	3866243.0	3931947.0	3996521.0	
14	56942108.0	56974100.0	57059007.0	57313203.0	
15	3499536.0	3470818.0	3443867.0	3415213.0	
16	436300.0	441525.0	446175.0	451630.0	
17	2367550.0	2337170.0	2310173.0	2287955.0	
18	390087.0	393028.0	395969.0	398582.0	
19	15925513.0	16046180.0	16148929.0	16225302.0	

	Population_2004	Population_2005	Population_2006	Population_2007	\
0	8171966.0	8227829.0	8268641.0	8295487.0	
1	10421137.0	10470617.0	10547958.0	10625780.0	
2	1018684.0	1037062.0	1055438.0	1073873.0	
3	10197101.0	10211216.0	10238905.0	10298828.0	
4	82516208.0	82469422.0	82376451.0	82266372.0	
5	5404623.0	5419432.0	5437272.0	5461438.0	
6	42923895.0	43653155.0	44397319.0	45226803.0	
7	1362550.0	1354775.0	1346810.0	1340680.0	
8	5228172.0	5266096.0	5266208.0	5280720.0	
9	62716306.0	63188395.0	63628261.0	64021377.0	
10	10955141.0	10987314.0	11020362.0	11048473.0	
11	4304600.0	4310145.0	4311159.0	4310217.0	
12	10107146.0	10007065.0	10071570.0	10055700.0	

```
import pandas as pd

file_path = '/content/drive/MyDrive/Colab Notebooks/recy_rates..csv'
recycling_data = pd.read_csv(file_path, header=7)

# Replace '.' (which likely indicates missing data) with NaN
recycling_data.replace('.', pd.NA, inplace=True)

# Mapping dictionary for country names to codes
country_code_mapping = {
    'Belgium': 'BEL',
    'Bulgaria': 'BGR',
    'Czechia': 'CZE',
    'Denmark': 'DNK',
    'Germany': 'DEU',
    'Estonia': 'EST',
    'Ireland': 'IRL',
    'Greece': 'GRC',
    'Spain': 'ESP',
    'France': 'FRA',
    'Croatia': 'HRV',
    'Italy': 'ITA',
    'Cyprus': 'CYP',
    'Latvia': 'LVA',
    'Lithuania': 'LTU',
    'Luxembourg': 'LUX',
    'Hungary': 'HUN',
    'Malta': 'MLT',
    'Netherlands': 'NLD',
    'Austria': 'AUT',
    'Poland': 'POL',
    'Portugal': 'PRT',
    'Romania': 'ROU',
    'Slovenia': 'SVN',
    'Slovakia': 'SVK',
    'Finland': 'FIN',
    'Sweden': 'SWE',
    # Add other countries here as needed
}

# Replace country names with codes
recycling_data['TIME'] = recycling_data['TIME'].map(country_code_mapping)

recycling_data = recycling_data.rename(columns={str(year): f'recy_rate_{year}' for year in range(2010, 2022)})

# Save cleaned recycling dataset to a CSV file
cleaned_file_path = '/content/drive/MyDrive/Colab Notebooks/cleaned_recycling_data.csv'
recycling_data.to_csv(cleaned_file_path, index=False)

# Print a message to confirm the file has been saved
print(f'Cleaned recycling data has been saved to: {cleaned_file_path}')

Cleaned recycling data has been saved to: /content/drive/MyDrive/Colab Notebooks/cleaned_recycling_data.csv
```



```
import pandas as pd

pop_file_path = '/content/drive/MyDrive/Colab Notebooks/cleaned_recy_data.csv'
cleaned_recycling_data = pd.read_csv(pop_file_path)

print(cleaned_recycling_data.head(20))
# Select columns for years 2010 to 2018
selected_columns = ['Year'] + [f'recy_rate_{year}' for year in range(2010, 2019)]
filtered_data = cleaned_recycling_data[selected_columns]

# Print filtered DataFrame to verify
print(filtered_data.head(20))

import pandas as pd

pop_file_path = '/content/drive/MyDrive/Colab Notebooks/cleaned_recy_data.csv'
cleaned_recycling_data = pd.read_csv(pop_file_path)

# Reset index to move country codes to a regular column
cleaned_recycling_data.reset_index(inplace=True)

# Rename 'Year' column to 'Country Code'
cleaned_recycling_data.rename(columns={'Year': 'Country Code'}, inplace=True)

# Select columns for years 2010 to 2018
selected_columns = ['Country Code'] + [f'recy_rate_{year}' for year in range(2010, 2019)]
filtered_data = cleaned_recycling_data[selected_columns]

# Define file path for previously combined dataset
file_path = '/content/drive/MyDrive/Colab Notebooks/combined_data.csv'

# Merge filtered data with combined data using 'Country Code' as common column
merged_data = pd.merge(combined_data, filtered_data, how='inner', on='Country Code')

# Specify the output file path to save merged dataset
#output_file_path = '/content/drive/MyDrive/Colab Notebooks/merged_combined_data.csv'

# Save merged dataset to a CSV file
merged_data.to_csv(output_file_path, index=False)
print(merged_data.head(5))
```

```
Country Name Country Code Indicator Name Indicator Code \
0 Austria AUT Population, total SP_POP.TOTL
1 Belgium BEL Population, total SP_POP.TOTL
2 Cyprus CYP Population, total SP_POP.TOTL
3 Czechia CZE Population, total SP_POP.TOTL
4 Germany DEU Population, total SP_POP.TOTL

Population_2000 Population_2001 Population_2002 Population_2003 \
0 8011566.0 8042293.0 8081957.0 8121423.0
1 10251250.0 10286570.0 10332785.0 10376133.0
2 948237.0 964830.0 982194.0 1000330.0
3 1018684.0 1037062.0 10190916.0 10193998.0
4 82211508.0 82349925.0 82488495.0 82534176.0

Population_2004 Population_2005 Population_2006 Population_2007 \
0 8171966.0 8227829.0 8268641.0 8295487.0
1 10421137.0 10478617.0 10547958.0 10625700.0
2 1018684.0 1037062.0 105438.0 1073873.0
3 10197101.0 10211216.0 10238905.0 1029828.0
4 82516260.0 82469422.0 82376451.0 82266372.0

Population_2008 Population_2009 Population_2010 Population_2011 \
0 8312496.0 8343329.0 8363486.0 8391643.0
1 10709973.0 10796493.0 10895586.0 11038264.0
2 1092390.0 1110974.0 1129686.0 1145086.0
3 10384003.0 10443936.0 10474410.0 10496088.0
4 82110097.0 81902307.0 81776930.0 80274983.0

Population_2012 Population_2013 Population_2014 Population_2015 \
0 8429991.0 8479823.0 8546356.0 8642699.0
1 11186932.0 11159407.0 11209057.0 11274196.0
2 1156556.0 1166968.0 1176995.0 1187280.0
3 10510785.0 10514272.0 10525347.0 10546059.0
4 80425823.0 80645605.0 80982500.0 81686611.0

Population_2016 Population_2017 Population_2018 Code GDP_2000 \
0 8736660.0 8797560.0 8840521.0 AUT 34796.258
1 11331422.0 11375158.0 11427054.0 BEL 33719.770
2 1197881.0 1208523.0 1218831.0 CYP 22326.799
3 10566332.0 10594438.0 10629928.0 CZE 17056.160
4 82348660.0 82657002.0 82905782.0 DEU 31367.205

GDP_2001 GDP_2002 GDP_2003 GDP_2004 GDP_2005 GDP_2006 \
0 33272.223 33823.586 36063.120 36957.113 37642.760 38866.840
1 33923.344 34419.695 34588.477 35740.350 36338.383 37051.902
2 23095.736 23828.072 24347.076 25355.240 26164.209 27142.018
3 17068.568 18431.027 19344.090 20555.035 22128.500 23888.164
4 34500.200 34590.930 34715.440 35320.715 36205.574 38014.137

GDP_2007 GDP_2008 GDP_2009 GDP_2010 GDP_2011 GDP_2012 GDP_2013 \
0 40305.273 40964.793 39463.656 40200.340 41446.0 42565.0 41375.0
1 38002.875 38117.348 36998.650 37739.330 38130.0 37906.0 37377.0
2 20109.996 20736.975 27735.033 27630.104 27272.0 26011.0 24519.0
3 25382.007 26186.045 25093.863 25922.395 26725.0 26474.0 26338.0
4 39752.207 40715.434 38962.938 41109.582 43109.0 43320.0 43413.0

GDP_2014 GDP_2015 GDP_2016 GDP_2017 GDP_2018 Plastic_Waste_2000 \
0 41338.0 41294.0 41445.0 42177.370 42906.070 26.210
```

```
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
import numpy as np

merged_data = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/merged_combined_data.csv')
```

```
# Calculate average GDP and population for each country over the years 2000 to 2018
gdp_columns = [f'GDP_{year}' for year in range(2000, 2019)]
pop_columns = [f'Population_{year}' for year in range(2000, 2019)]
```

```
merged_data['Average_GDP'] = merged_data[gdp_columns].mean(axis=1)
merged_data['Average_Population'] = merged_data[pop_columns].mean(axis=1)
```

```
# Use averages for clustering
data_to_cluster = merged_data[['Average_GDP', 'Average_Population']]
```

```
# Standardize data
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data_to_cluster)
```

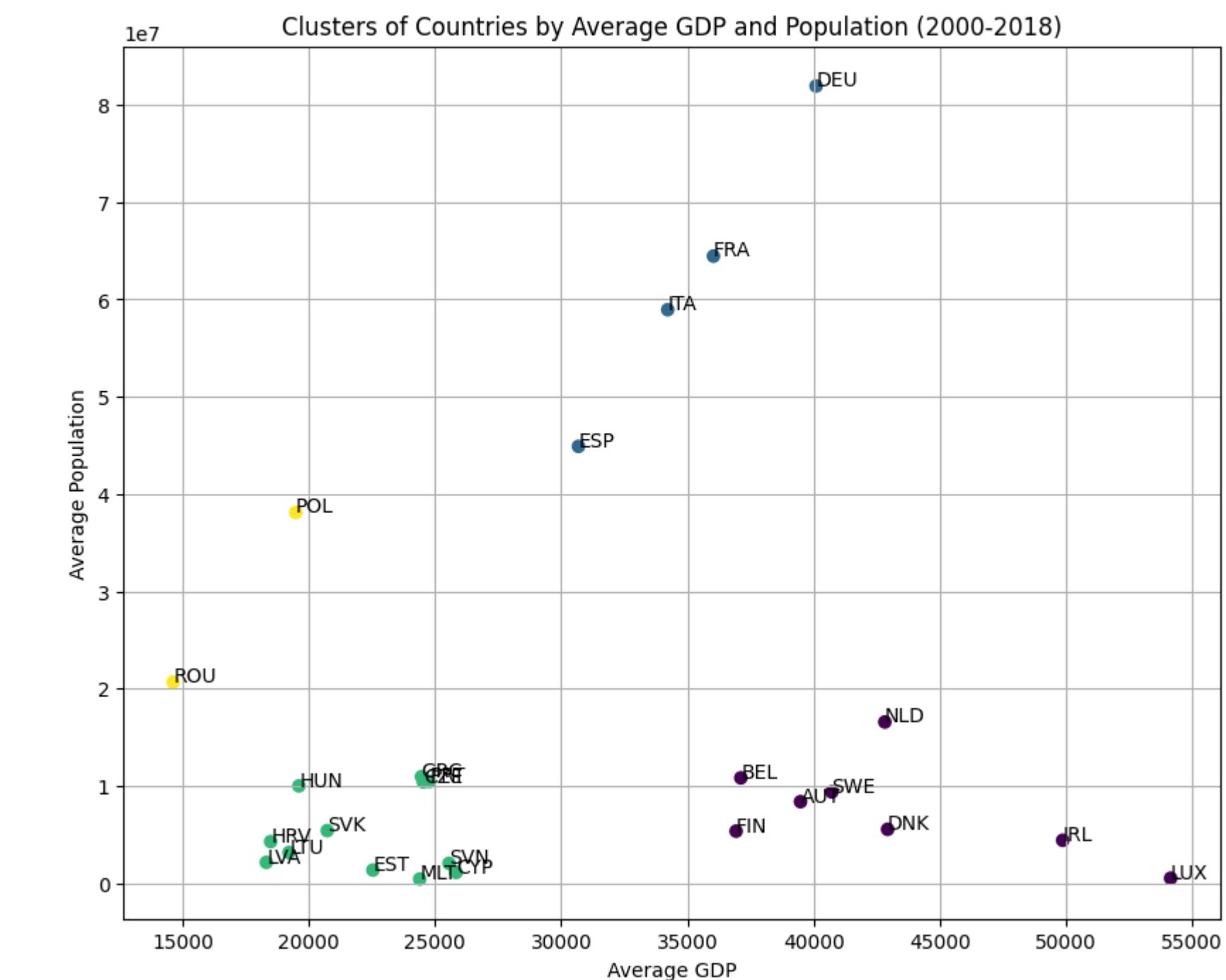
```
# Apply K-Means clustering
kmeans = KMeans(n_clusters=4, random_state=0)
clusters = kmeans.fit_predict(data_scaled)
```

```
# Create a scatter plot
plt.figure(figsize=(10, 8))
plt.scatter(data_to_cluster['Average_GDP'], data_to_cluster['Average_Population'], c=clusters, cmap='viridis')
```

```
# Annotate country codes
for i, txt in enumerate(merged_data['Country Code']):
    plt.annotate(txt, (data_to_cluster['Average_GDP'][i], data_to_cluster['Average_Population'][i]))
```

```
plt.title('Clusters of Countries by Average GDP and Population (2000-2018)')
plt.xlabel('Average GDP')
plt.ylabel('Average Population')
plt.grid(True)
plt.show()
```

```
/usr/local/lib/python3.10/dist-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()
```



```
plastic_waste_columns = [f'Plastic_Waste_{year}' for year in range(2000, 2019)]
```

```
# Check if all columns are present
if all(column in merged_data.columns for column in plastic_waste_columns):
    # Calculate the average plastic waste and population for each country over the years 2000 to 2018
    merged_data['Average_Plastic_Waste'] = merged_data[plastic_waste_columns].mean(axis=1)
```

```
# Select average data for clustering
data_to_cluster = merged_data[['Average_Plastic_Waste', 'Average_Population']].copy()
```

```
# Standardize data
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data_to_cluster)
```

```
# Apply K-Means clustering
kmeans = KMeans(n_clusters=4, random_state=0)
clusters = kmeans.fit_predict(data_scaled)
```

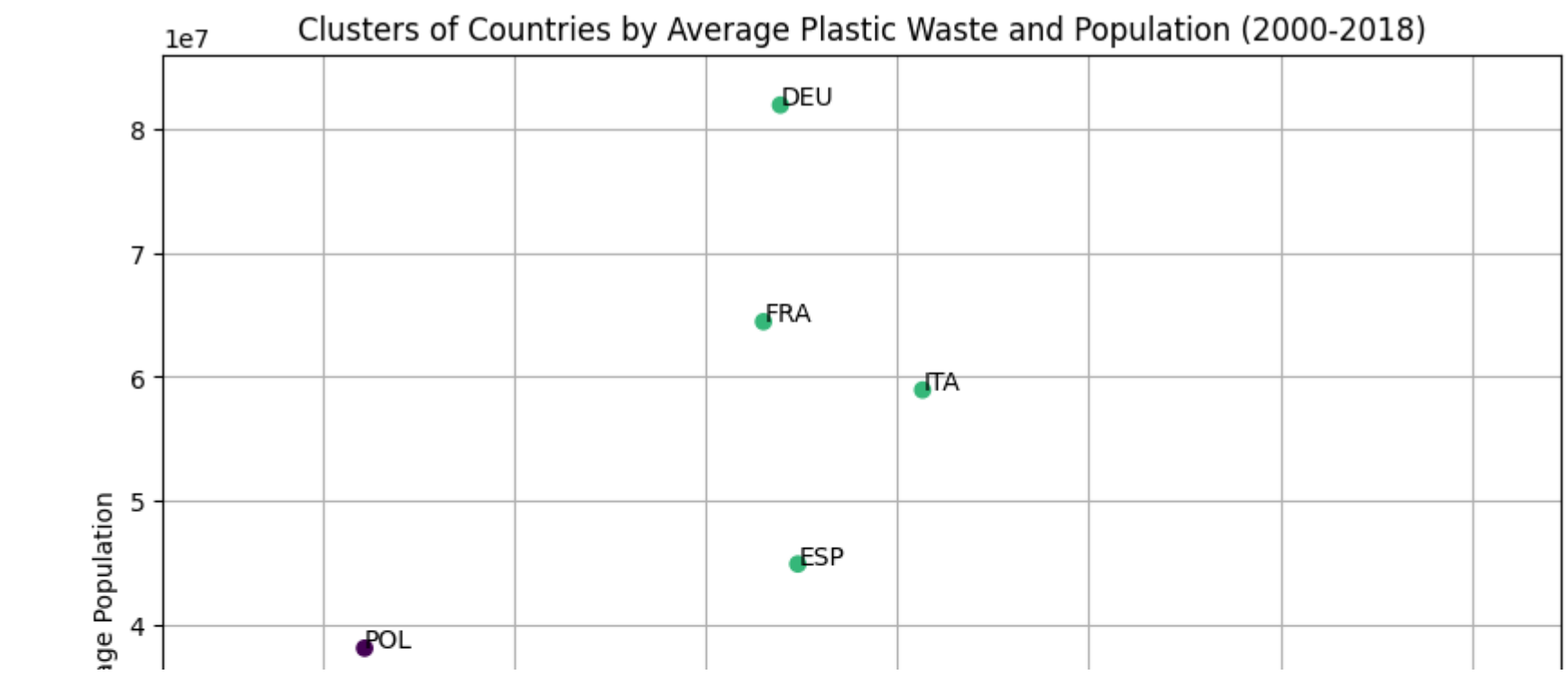
```
# Create a scatter plot
plt.figure(figsize=(10, 8))
plt.scatter(data_to_cluster['Average_Plastic_Waste'], data_to_cluster['Average_Population'], c=clusters, cmap='viridis')
```

```
# Annotate the country codes
for i, txt in enumerate(merged_data['Country Code']):
    plt.annotate(txt, (data_to_cluster['Average_Plastic_Waste'][i], data_to_cluster['Average_Population'][i]))
```

```
plt.title('Clusters of Countries by Average Plastic Waste and Population (2000-2018)')
plt.xlabel('Average Plastic Waste')
plt.ylabel('Average Population')
plt.grid(True)
plt.show()
```

```
else:
    print("One or more columns for plastic waste are missing from the DataFrame.")
```

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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



```
# Calculate average GDP, population, and plastic waste for each country
gdp_columns = [f'GDP_{year}' for year in range(2000, 2019)]
pop_columns = [f'Population_{year}' for year in range(2000, 2019)]
plastic_waste_columns = [f'Plastic_Waste_{year}' for year in range(2000, 2019)]

# Check if all columns are present
if all(column in merged_data.columns for column in gdp_columns + pop_columns + plastic_waste_columns):
    merged_data['Average_GDP'] = merged_data[gdp_columns].mean(axis=1)
    merged_data['Average_Population'] = merged_data[pop_columns].mean(axis=1)
    merged_data['Average_Plastic_Waste'] = merged_data[plastic_waste_columns].mean(axis=1)

# Select average data for correlation
data_for_correlation = merged_data[['Average_GDP', 'Average_Population', 'Average_Plastic_Waste']]

# Calculate Pearson correlation matrix
correlation_matrix = data_for_correlation.corr()
print(correlation_matrix)
```

```
else:
    raise KeyError("One or more columns for GDP, population, or plastic waste are missing from the DataFrame.")
```

	Average_GDP	Average_Population	Average_Plastic_Waste
Average_GDP	1.000000	0.145851	0.700062
Average_Population	0.145851	1.000000	0.183462
Average_Plastic_Waste	0.700062	0.183462	1.000000

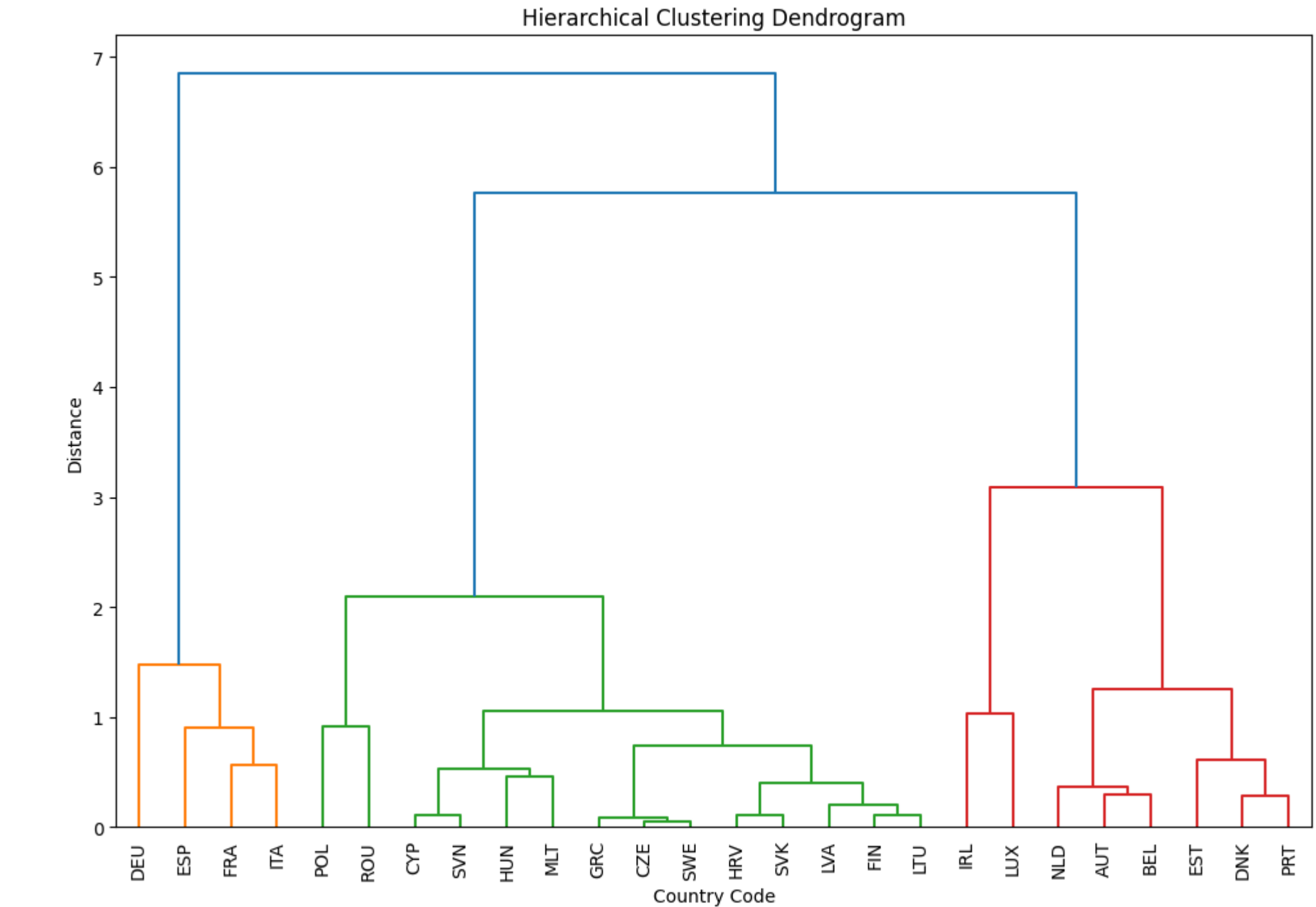
```
from scipy.cluster.hierarchy import dendrogram, linkage

# We already have scaled data in 'data_scaled' from previous standardization step.
# Generate linkage matrix for hierarchical clustering
Z = linkage(data_scaled, method='ward')

# Set up matplotlib figure
plt.figure(figsize=(12, 8))

# Generate and plot dendrogram
dendrogram(Z, labels=merged_data['Country Code'].values, leaf_rotation=90, leaf_font_size=10)

plt.title('Hierarchical Clustering Dendrogram')
plt.xlabel('Country Code')
plt.ylabel('Distance')
plt.show()
```



```
# Re-apply K-Means clustering to PCA results for color coding
kmeans = KMeans(n_clusters=4, random_state=0)
pca_clusters = kmeans.fit_predict(pca_result)
pca_df['Cluster'] = pca_clusters # Add cluster assignment to PCA results DataFrame

# Plot PCA results with color coding by cluster
plt.figure(figsize=(12, 10))
scatter = plt.scatter(pca_df['PC1'], pca_df['PC2'], c=pca_df['Cluster'], cmap='viridis')

# Add annotations for each point
for i, txt in enumerate(pca_df['Country Code']):
    plt.annotate(txt, (pca_df['PC1'][i], pca_df['PC2'][i]), fontsize=9)

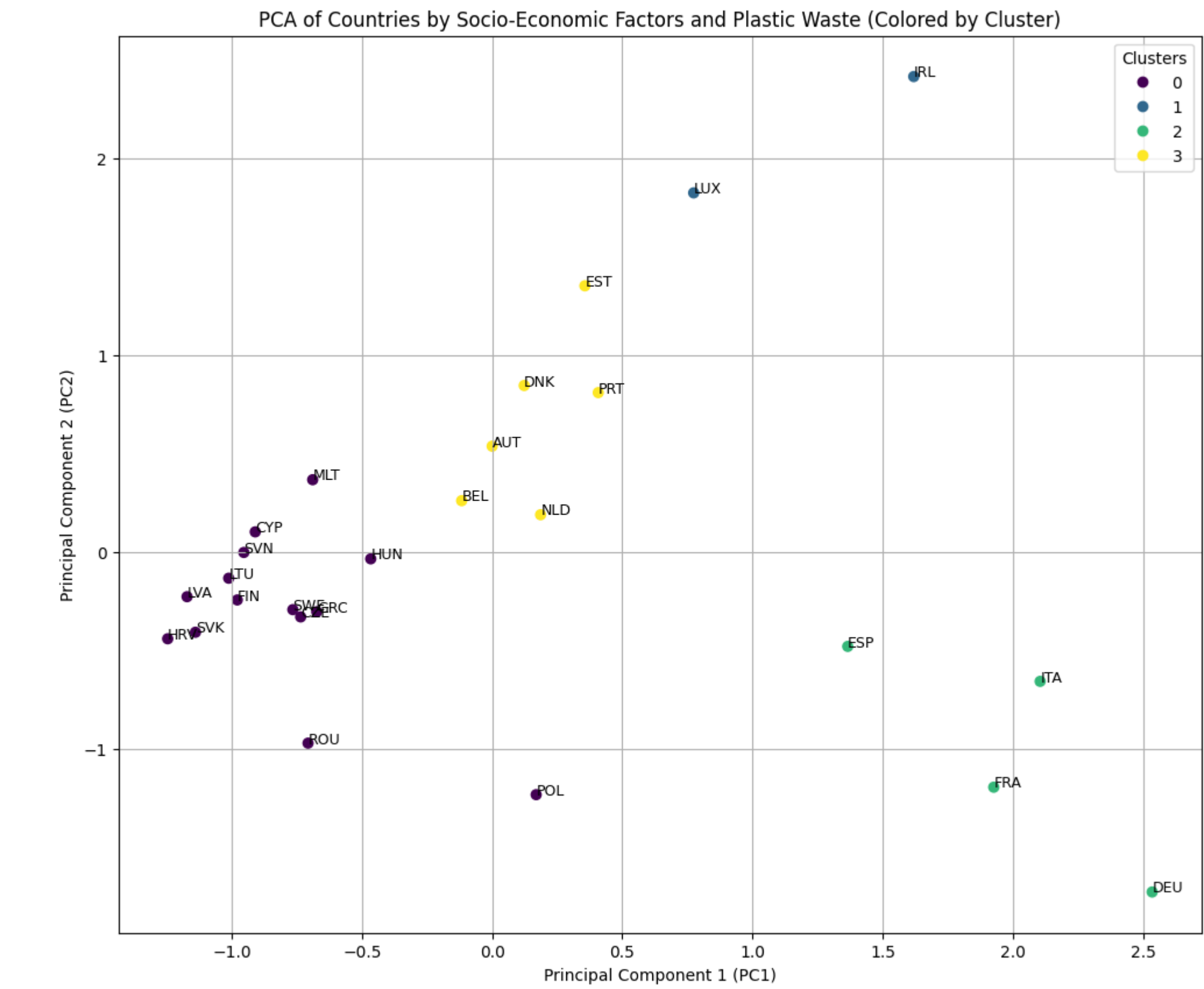
# Add legend for clusters
plt.legend(*scatter.legend_elements(), title="Clusters")

# Set title and axis labels
plt.title('PCA of Countries by Socio-Economic Factors and Plastic Waste (Colored by Cluster)')
plt.xlabel('Principal Component 1 (PC1)')
plt.ylabel('Principal Component 2 (PC2)')
plt.grid(True)

plt.show()
```

```
# Print explained variance
explained_variance = pca.explained_variance_ratio_
explained_variance
```

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:878: 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



array([0.59173111, 0.48826889])