# Assignment 3 - World Bank Data Analysis

In [10]:

```python
"""Importing all neceassary libraries"""
import wbgapi as wb
import pandas as pd
import seaborn as sns
import numpy as np
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn import linear_model
from scipy import spatial
import itertools as iter


def load_clean_data():
    """ Loading and clearning the data """

    dataframe = pd.read_csv("WDIData.csv")
    del dataframe["Unnamed: 66"]
    dataframe = dataframe.fillna(0)

    return dataframe

def visualize_emu_gpd():
    """ EMU countries gdp visualization """

    gdppercap = wb.data.DataFrame('NY.GDP.PCAP.CD', wb.region.m
    g5 = gdppercap.sort_values(by=['YR2020'], ascending = False
    ax = gdppercap.T.plot(color = 'lightgray', legend=False)
    g5.T.plot(ax=ax, figsize=(15,5))

def get_visualization(df):
    """This Function will visualized data according to the GDP

    df = df.head(7000)
    df = df.fillna(0)

    palette = sns.color_palette("Paired", 10)
    sns.set_palette(palette)

    #we take only data, not additional informations
    df = df[0:-5]
    df.replace('..', np.nan, inplace=True)

    col_list = df.columns[4:].values
```

```python
48        df[col_list]=df[col_list].apply(pd.to_numeric)
49        #reindex all table, create pivot view
50        pv2 = pd.pivot_table(df,index=['Indicator Name','Country Co
51        # set the years
52        pv2.columns = np.arange(1960,2022)
53        palette = sns.color_palette("Paired", 10)
54        sns.set_palette(palette)
55
56        pv2.loc['GDP (current US$)'].T.plot(alpha=1, rot=45)
57        pv2.loc['GDP per capita (current US$)'].T.plot(alpha=0.8, ᵣ
58        pv2.loc['GDP per capita (current US$)'].T.plot(alpha=0.75,
59        pv2.loc['GDP growth (annual %)'].T.plot(alpha=0.75, rot=45)
60
61
62    def gdp_clustering(dataframe):
63        """ GDP clustering using K-Means """
64
65        years = dataframe.columns[4:].tolist()
66        new_data = dataframe.copy()
67        new_data = new_data[:10000]
68        year_values = []
69        for index, row in new_data.iterrows():
70            one_row_gdp = []
71            if "GDP" in row["Indicator Name"]:
72                if row.any():
73                    for year in years:
74                        one_row_gdp.append(row[year])
75                    one_row_gdp.append(1)
76                    year_values.append(np.array(one_row_gdp))
77
78        year_values = np.array(year_values)
79
80        kmeans = KMeans(n_clusters=2, random_state=0)
81        clusters = kmeans.fit_predict(year_values)
82        kmeans.cluster_centers_.shape
83        return clusters[:1000]
84
85
86    def co2_clustering(dataframe):
87        """ CO2 clustering using K-Means """
88
89        years = dataframe.columns[4:].tolist()
90        new_data = dataframe.copy()
91        new_data = new_data[:10000]
92        year_values = []
93        for index, row in new_data.iterrows():
94            one_row_gdp = []
95            if "CO2" in row["Indicator Name"]:
96                if row.any():
97                    for year in years:
98                        one_row_gdp.append(row[year])
99                    one_row_gdp.append(1)
100            year_values.append(np.array(one_row_gdp))
101
```

```python
102        year_values = np.array(year_values)
103
104        kmeans = KMeans(n_clusters=2, random_state=0)
105        clusters = kmeans.fit_predict(year_values)
106        kmeans.cluster_centers_.shape
107
108        return clusters[:1000]
109
110
111    def normalize_values(col):
112        """ Min Max normalization """
113
114        max_value = col.max()
115        min_value = col.min()
116        new_col = (col - min_value) / (max_value - min_value)
117        return new_col
118
119
120    def prediction(dataframe, years):
121        """ Linear Regression grouped by years """
122
123        X = np.array([normalize_values(dataframe[year]) for year in
124        y = years
125
126        X_train, X_test, y_train, y_test = train_test_split(X, y, t
127
128        model = linear_model.LinearRegression()
129        model = model.fit(X_train, y_train)
130        predicted_data = model.predict(X_test)
131        predicted_data = np.round_(predicted_data)
132
133        MSE = mean_squared_error(y_test,predicted_data)
134        PD = predicted_data
135        return MSE
136
137
138    def country_clustering(dataframe):
139        """ Country Grouping and Feature Vectors """
140
141        years = dataframe.columns[4:].tolist()
142        new_data = dataframe.copy()
143        new_data = new_data[-20000:]
144        countries = {}
145
146        for index, row in new_data.iterrows():
147            one_row_country = []
148            if row.any():
149                for year in years:
150                    one_row_country.append(row[year])
151
152            if row['Country Name'] in countries:
153                countries[row['Country Name']].extend(one_row_count
154            else:
155                countries[row['Country Name']] = [one_row_country]
```

```python
156
157          return countries
158
159
160      def cosine(country1, country2):
161          """ Cosine Similarity Function """
162
163          result = 1 - spatial.distance.cosine(country1, country2)
164
165          return result
166
167      def err_ranges(x, func, param, sigma):
168          """
169          Calculates the upper and lower limits for the function, pa
170          sigmas for single value or array x. Functions values are c
171          all combinations of +/- sigma and the minimum and maximum
172          Can be used for all number of parameters and sigmas >=1.
173
174          This routine can be used in assignment programs.
175          """
176
177          # initiate arrays for lower and upper limits
178          lower = func(x, *param)
179          upper = lower
180
181          uplow = []    # list to hold upper and lower limits for para
182          for p,s in zip(param, sigma):
183              pmin = p - s
184              pmax = p + s
185              uplow.append((pmin, pmax))
186
187          pmix = list(iter.product(*uplow))
188
189          for p in pmix:
190              y = func(x, *p)
191              lower = np.minimum(lower, y)
192              upper = np.maximum(upper, y)
193
194          return lower, upper
195
```
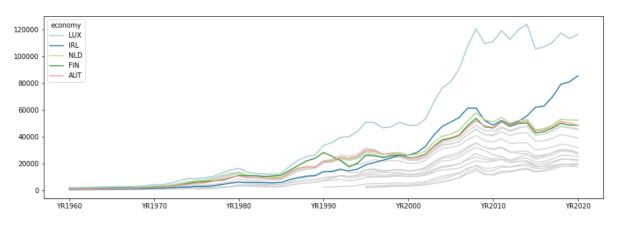
In [11]:

```python
if __name__ == "__main__":
    dataframe = load_clean_data()
    gdp_df = gdp_clustering(dataframe)
    print(f"GDP Clusters: {gdp_df}")
    print("\n")
    co2_df = co2_clustering(dataframe)
    print(f"CO2 Clusters: {co2_df}")
    print("\n")
    years = dataframe.columns[4:]
    error = prediction(dataframe, years)
    print(f"Prediction Error: {error}")
    print("\n")

    countries = country_clustering(dataframe)
    ukrain = np.array(countries["Ukraine"][-1000:])
    vietnam = np.array(countries["Vietnam"][-1000:])
    similarity = cosine(ukrain, vietnam)
    print(f"Cosine Similarity of Ukraine and Vietnam: {similari
    print("\n")
    print("Visualizing GDP Growth with 10 years span")
    visualize_emu_gpd()
    print("Visualizing GDP per capita")
    get_visualization(dataframe)
```

```
GDP Clusters: [1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
```

```
1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1 1
1 1 1 1
 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1
 1 1 1 1 1 1]
```

```
CO2 Clusters: [0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0
 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0
0 0 0 0
 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0]
```
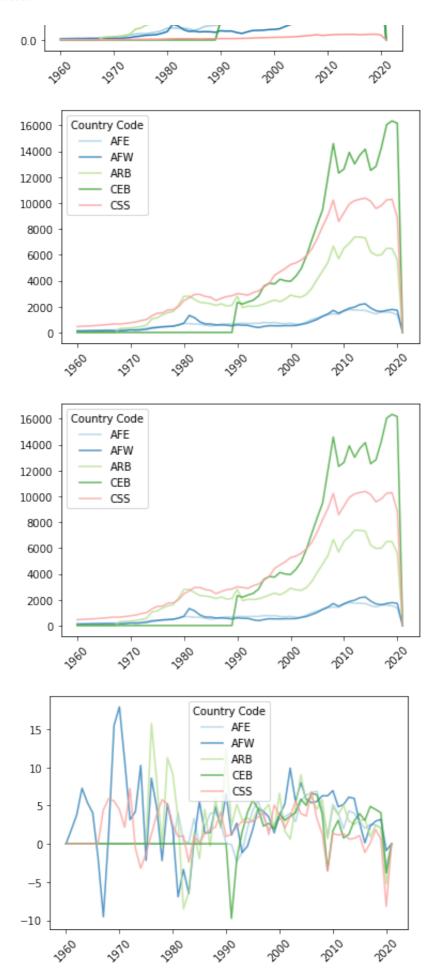
Prediction Error: 3.8125

Cosine Similarity of Ukraine and Vietnam: 0.9262062550042696

Visualizing GDP Growth with 10 years span
Visualizing GDP per capita

In [ ]: