

# Assignment 1

## Students

- Total number of hours spent on assignment per person: 7 hours.
- Davíð Freyr Björnsson
- Eric Guldbrand

## 1

Write a Python program that draws a scatter plot of GDP per capita vs life expectancy. State any assumptions and in combining data. Answer these questions:

1. Which countries have a life expectancy higher than one standard deviation above the mean?
2. Which countries have high life expectancy but have low GDP?
3. Does every strong economy have high life expectancy?

## Description of the data

- Data for GDP per capita comes from the World Bank and contains the period 1990 to 2017.
- Data for life expectancy comes from OWID. OWID has combined two data sources: Clio-Infra's dataset for Population Division for data from 1950 to 2015.

```
1 # Load Data
2
3 import pandas as pd
4
5 gdpPerCapitaURL = 'https://raw.githubusercontent.com/ecen/data-science-intro/master/gdpPerCapita.csv'
6 gdpPerCapitaDF = pd.read_csv(gdpPerCapitaURL)
7
8 lifeExpectancyURL = 'https://raw.githubusercontent.com/ecen/data-science-intro/master/lifeExpectancy.csv'
9 lifeExpectancyDF = pd.read_csv(lifeExpectancyURL)

1 # Look at the structure of the data
2 gdpPerCapitaDF.describe()
```



```
1 # Rename the GDP per capita
2 gdpPerCapitaDF.rename(columns={"GDP per capita (int.-$) (constant 2011 international" data-bbox="94 935 954 964"/>
```

```

3 lifeExpectancyDF.rename(columns={"Life expectancy (Clio-Infra up to 1949; UN Populatio
1 # Select and combine data
2
3 gdpPerCapita2015 = gdpPerCapitaDF.loc[gdpPerCapitaDF['Year'] == 2015]
4 lifeExpectancy2015 = lifeExpectancyDF.loc[lifeExpectancyDF['Year'] == 2015]
5
6 gdpVsLife = pd.merge(gdpPerCapita2015, lifeExpectancy2015)

```

It was decided to use data from 2015, the most recent year that existed in both the selected datasets.

The 2015 data from both datasets was then merged, retaining only countries that existed in both.

```

1 # Create a scatterplot, with GDP per capita in 2015 on the x-axis
2 # and life expectancy in 2015 on the y-axis
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 plt.scatter(gdpVsLife['gdpPerCapita'], gdpVsLife['lifeExp'])
7 plt.xlabel("GDP per capita (2015)")
8 plt.ylabel("life expectancy in 2015")
9 plt.show()

```



## ▼ a. Which countries have a life expectancy higher than one standard deviation

```

1 # Selected countries with a life expectancy of one standard deviation or more over the
2 lifeStd = gdpVsLife['lifeExp'].std()
3 lifeMean = gdpVsLife['lifeExp'].mean()
4
5 lifeAboveOneStd = gdpVsLife.loc[gdpVsLife['lifeExp'] > lifeMean + lifeStd]
6
7 lifeAboveOneStd.Entity

```



## ▼ b. Which countries have high life expectancy but have low GDP?

```
1 # Select countries with high life expectancy and low gdp
2 # We define high and low as the 75th respective 25th percentile
3
4 highGdp75 = gdpVsLife.gdpPerCapita.quantile(0.75)
5 highGdp50 = gdpVsLife.gdpPerCapita.quantile(0.5)
6 lowGdp50 = gdpVsLife.gdpPerCapita.quantile(0.5)
7 lowGdp25 = gdpVsLife.gdpPerCapita.quantile(0.25)
8 highLife75 = gdpVsLife.lifeExp.quantile(0.75)
9 highLife50 = gdpVsLife.lifeExp.quantile(0.5)
10
11 lifeHighGdpLow = gdpVsLife.loc[(gdpVsLife['lifeExp'] > highLife75) & (gdpVsLife['gdpP
12
13 lifeHighGdpLow.Entity
```



Given a more strict definition of high and low, 75th and 25th percentiles respectively, there are no countries with l

```
1 # Select countries with high life expectancy and low gdp
2 # We define high and low as the 50th respective 50th percentile
3
4 lifeHighGdpLow = gdpVsLife.loc[(gdpVsLife['lifeExp'] > highLife50) & (gdpVsLife['gdpP
5
6 lifeHighGdpLow.Entity
```



Given a less strict definition, with high and low being considered above and below the median (50th percentile) life expectancy but low GDP.

### ▼ c. Does every strong economy have high life expectancy?

```
1 gdpHigh = gdpVsLife.loc[gdpVsLife['gdpPerCapita'] > highGdp75]
2 lifeHighGdpHigh = gdpVsLife.loc[(gdpVsLife['lifeExp'] > highLife75) & (gdpVsLife['gdpPerCapita'] > highGdp75)]
3
4 gdpHigh.shape[0] - lifeHighGdpHigh.shape[0]
```



```
1 plt.scatter(gdpVsLife['gdpPerCapita'], gdpVsLife['lifeExp'])
2 plt.xlabel("GDP per capita (2015)")
3 plt.ylabel("life expectancy in 2015")
4 plt.axvline(highGdp75, color='orange')
5 plt.axhline(highLife75, color='orange')
6 plt.text(highGdp75 - 10000, 87, '75th percentile')
7 plt.text(130000, highLife75, '75th percentile')
8 plt.show()
```



```
1 ## Does the result change much if we alter our definition of low and high?
2 ## Set high and low both as the 50th percentile
3
4 gdpHigh = gdpVsLife.loc[gdpVsLife['gdpPerCapita'] > highGdp50]
5
6 lifeHighGdpHigh = gdpVsLife.loc[(gdpVsLife['lifeExp'] > highLife50) & (gdpVsLife['gdpPerCapita'] > highGdp50)]
7
8 gdpHigh.shape[0] - lifeHighGdpHigh.shape[0]
```



```
1 plt.scatter(gdpVsLife['gdpPerCapita'], gdpVsLife['lifeExp'])
2 plt.xlabel("GDP per capita (2015)")
3 plt.ylabel("life expectancy in 2015")
4 plt.axvline(highGdp50, color='orange')
```

```

5 plt.axhline(highLife50, color='orange')
6 plt.text(highGdp50 - 10000, 87, '50th percentile')
7 plt.text(130000, highLife50, '50th percentile')
8 plt.show()

```



Not every strong economy has high life expectancy. If we define high as above 75th percentile and low as below 75th percentile, we can see that some strong economies and not high life expectancy. If we define high and low to be above and below the median then

## 2. Economic growth compared to tertiary education enr

```

1 # Annual growth of GDP per capita, adjusted for inflation, 1961 to 2014: https://ourw
2 gdpPerCapitaGrowth = pd.read_csv('https://raw.githubusercontent.com/ecen/data-science
3 gdpPerCapitaGrowth.rename(columns={"GDP per capita growth (annual %) (%)": "growth"},
4
5 # Tertiary enrollment: Total enrollment in tertiary education (ISCED 5 to 8), regardl
6 # expressed as a percentage of the total population of the five-year age group follow
7 # Years: 1970 to 2015. Source: https://data.worldbank.org/data-catalog/ed-stats, UNES
8 # expressed as a percentage of the total population of the five-year age group follow
9 enrollment = pd.read_csv('https://raw.githubusercontent.com/ecen/data-science-intro/m
10 enrollment.rename(columns={"Gross enrolment ratio, tertiary, both sexes (%) (%)": "enr

```

```

1 def plotYears(startYear, endYear):
2     growthData = gdpPerCapitaGrowth[(gdpPerCapitaGrowth['Year'] >= startYear) & (gdpPer
3     enrollmentData = enrollment[(enrollment['Year'] >= startYear) & (enrollment['Year']
4     data = growthData.merge(enrollmentData, on='Entity')
5
6     # Remove outliers to make it easier to compare the distribution of the majority
7     axes = plt.gca()
8     axes.set_xlim([-15, 15])
9     axes.set_ylim([0, 80])
10    plt.scatter(data.growth, data.enrollment)
11
12    plt.xlabel("Average growth in GDP per Capita (%)")
13    plt.ylabel("Average growth in tertiary enrollment (%)")
14    plt.title("Average growth in GDP per capita" + "\n" + "and tertiary enrollment for
15              + str(startYear) + " - " + str(endYear))
16    plt.show()
17
18 plotYears(1970, 1980)
19 plotYears(1980, 1990)
20 plotYears(1990, 2000)
21 plotYears(2000, 2010)

```



Figures show the 10-year average growth in GDP per capita in percent versus tertiary enrollment in percent in the school leaving.

A few outliers have been excluded from these figures to ease comparison of the distributions between the 10-year periods. Over time we see an increase in tertiary education enrollment (again, of those who have recently finished secondary education). The first 10 year period (1970-1980) seems to indicate two groups. Countries with low tertiary enrollment have a wider spread, while countries with higher tertiary enrollment all have a more moderate spread. In later 10 year periods, this split seems to disappear.