

# Relation between Gross Domestic Product (GDP) per capita and life satisfaction of different countries

**Damian Ejlli**

## 1 Introduction

Happiness is probably the most important thing that many people aim to achieve in life but what are the factors that contribute to it is not immediately known. In fact, happiness very likely depends on several factors that contribute to it and some of these factors are more important than the others. Obviously happiness is very subjective and there are many ways to quantify it and one of the most important quantifiers of happiness is the perceived life satisfaction of a group of people. In this article I want to study what is the relation between the Gross Domestic Product (GPD) and the life satisfaction of a given country.

## 2 Data file loading and overview

To study the relation between GPD and life satisfaction, first is needed to have the tabular data of the GPD and life satisfaction of different countries and second it is necessary to prepare the data for the analysis. The tabular data of the GPD and life satisfaction can be easily found on internet respectively in the (IMF website, GPD per capita in US dollars) and (OECD website, Better Life Index). Both these tabular data files are provided with the Python code for analysis. Before starting the analysis, it is very important to have a look at the data present in the files in the tabular form by using for example Google Sheets. In Fig. 1, a section of the Better Life Index aggregate data on Google Sheet is shown. The file contains different information regarding Better Life Index but I am mostly interested in some of these data. I am interest in the column A that gives the country name, column D that gives the indicator of better life which in our case is (Life satisfaction), column G that gives inequality measurement and column O that gives the numeric value of each better life index present in column D.

In Fig. 2, a section of the Better Life Index is shown where in the column D we have the desired category of index, namely that of Life satisfaction index. The first data appears in the row nr. 1815 with Australia being the first country. If we roll down little with the file, we can see that there are data of the Life satisfaction for men, women and total (women+men). Here I work only with the total data which have a corresponding average score (columns I and J). These data extend from row nr. 1815 to 1853 and will be the data that I will use in what follows. The Life satisfaction data are displayed in the score scale from 0 to 10 where the value of 10 means extremely satisfied while the score value of 0 means extremely unsatisfied.

LOCATION																
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	LOCATION	Country	INDICATOR	Indicator	MEASURE	Measure	INEQUALITY	Inequality	Unit Code	Unit	PowerCode Cod	PowerCode	Reference Period	Reference Period	Value	Flag Codes
2	AUS	Australia	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					5.4	
3	AUT	Austria	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.5	
4	BEL	Belgium	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.7	
5	CAN	Canada	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					6	
6	CZE	Czech Republic	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.1	
7	DNK	Denmark	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					4.2	
8	FIN	Finland	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.9	
9	FRA	France	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					7.6	
10	DEU	Germany	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					2.7	
11	GRC	Greece	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					29.8	
12	HUN	Hungary	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					4.7	
13	ISL	Iceland	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					0.7	
14	IRL	Ireland	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					7.8	
15	ITA	Italy	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					12.3	
16	JPN	Japan	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					1.4	
17	KOR	Korea	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					2.6	
18	LUX	Luxembourg	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					1.7	
19	MEX	Mexico	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					5.5	
20	NLD	Netherlands	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					4.8	
21	NZL	New Zealand	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					4.7	
22	POL	Poland	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					5.7	
23	PRT	Portugal	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					10	
24	SVK	Slovak Republic	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					9.9	
25	ESP	Spain	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					23.1	
26	SWE	Sweden	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.2	
27	TUR	Turkey	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					12.5	
28	GBR	United Kingdom	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					4.5	
29	USA	United States	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					7.7	
30	CHL	Chile	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					8.7	
31	EST	Estonia	JE_LMIS	Labour market in L	Value	TOT	Total	PC	Percentage	0 Units					3.8	

**Figure 1:** Appearance of the first 31 rows of OECD Better Life Index aggregate data on Google Sheets. Here I am interested in the column A data (Country), column D data (Indicator), column G data (INEQUALITY) and column O data (Value).

LOCATION																
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1812	LVA	Latvia	HS_SFRH	Self-reported h	L	Value	LW	Low	PC	Percentage	0 Units					28
1813	SVN	Slovenia	HS_SFRH	Self-reported h	L	Value	LW	Low	PC	Percentage	0 Units					53
1814	OECD	OECD - Total	HS_SFRH	Self-reported h	L	Value	LW	Low	PC	Percentage	0 Units					61
1815	AUS	Australia	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.3	
1816	AUT	Austria	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.1	
1817	BEL	Belgium	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.9	
1818	CAN	Canada	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.4	
1819	CZE	Czech Republic	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.7	
1820	DNK	Denmark	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.6	
1821	FIN	Finland	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.6	
1822	FRA	France	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.5	
1823	DEU	Germany	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7	
1824	GRC	Greece	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.4	
1825	HUN	Hungary	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.6	
1826	ISL	Iceland	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.5	
1827	IRL	Ireland	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7	
1828	ITA	Italy	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6	
1829	JPN	Japan	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.9	
1830	KOR	Korea	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.9	
1831	LUX	Luxembourg	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.9	
1832	MEX	Mexico	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.5	
1833	NLD	Netherlands	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.4	
1834	NZL	New Zealand	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.3	
1835	NOR	Norway	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.6	
1836	POL	Poland	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.1	
1837	PRT	Portugal	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.4	
1838	SVK	Slovak Republic	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.2	
1839	ESP	Spain	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.3	
1840	SWE	Sweden	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.3	
1841	CHE	Switzerland	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				7.5	
1842	TUR	Turkey	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				5.5	
1843	GBR	United Kingdom	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.8	
1844	USA	United States	SW_LIFS	Life satisfaction	L	Value	TOT	Total	AVSCORE	Average score	0 Units				6.9	

**Figure 2:** Google sheet section containing the (Life satisfaction) indicator in column D of the Better Life Index aggregate data file.

**Figure 3:** Google sheet section containing the GDP per capita in (US dollars) for different countries and years.

**Figure 4:** Google sheet section containing the GDP per capita in (US dollars) where I select as a matter of example the column corresponding to the year 2015 for my analysis.

In Fig. 3 a section of the data file containing the information of the GPD per capita is (US dollars) for different countries and years in Google sheets is shown. We can see that the data appear quite messy and they overlap with each other. The data sheet contains different data and I will be interested only in some of them. The heading of the GPD per capita file starts at the row nr. 1 in Fig. 3 “Data Source” and “World Development Indicators” appear. In what follows, I will be interested in the data with heading in row nr. 5 which among them are the “Country Name”, “Indicator Name” and years (in numerical values). In Fig. 4 a section of the GPD per capita file corresponding to the years from 2006 to 2019 is shown for different countries. As a matter of example, in my analysis I choose<sup>1</sup> the year 2015 which is the selected column in blue in Fig. 4.

### 3 Data wrangling

In sec. 2, I showed some sections of the Better Life Index and GPD per capita files as they appear when visualized with Google Sheets. In this section, I show how to clean and prepare the data for the analysis by using Python 3. It is important to set since know the goal of my analysis and the way to reach that goal. The goal of my analysis is that to collect and tabulate the data of the life satisfaction index and GDP per capita for each country name and find the relationship (if there is one) between the data.

I start first by importing the libraries and/or modules that I use in my python analysis by using the “Jupyter Notebook”:

- `import numpy as np`
- `import matplotlib.pyplot as plt`
- `import pandas as pd`
- `import scipy as sp`
- `import sklearn.neighbors`

Here I have imported the NumPy, Matplotlib, Pandas, Scipy and Sklearn modules/libraries and have used the standard abbreviations to call them. As I discussed in sec. 2 both data files contain many data categories that are not needed for this analysis and I select only those data that are necessary to reach my goal. I start first by loading the files in Jupyter Notebook by using Pandas library

- `GDP=pd.read_csv(“gdp per capita.csv”, delimiter=“,”, header=2)`
- `LS=pd.read_csv(“better life index.csv”)`

The next step is that of data wrangling by using the Pandas library. In the GDP per capita file I set the “Country name” as the index in the Pandas dataframe and the value column being the “ GDP per capita 2015 (USD)”. Here I use Pandas to rename some columns and set as dataframe index the column that I am interested in of the original tabular data. Details of these operation are presented in the accompanying Jupyter Notebook file. I use similar data wrangling for the Better Life Index file where

---

<sup>1</sup>One is free to choose another year for the analysis if wishes so.

I rename some columns and set the country column as index of the Pandas dataframe and column value being the “Life Satisfaction Value”. In Fig. 5 sections of the Pandas dataframes are shown. On the left a section of the “GDP per capita 2015 (USD)” dataframe is shown and on the right only the first ten entries of the “Life Satisfaction Value” dataframe is shown. After data wrangling the “GDP per capita 2015 (USD)” dataframe has 264 country entry values of the GDP per capita as can be seen from the left dataframe in Fig. 5 while the “Life Satisfaction Value” dataframe has in total 40 country value entries of the life satisfaction values.

GPD per capita 2015 (USD)		Life Satisfaction Value	
Country Name		Country Name	
Aruba	27980.880695	Australia	7.3
Afghanistan	578.466353	Austria	7.1
Angola	4166.979684	Belgium	6.9
Albania	3952.801215	Canada	7.4
Andorra	35762.523074	Czech Republic	6.7
...	...	Denmark	7.6
Kosovo	3603.025501	Finland	7.6
Yemen, Rep.	1602.037841	France	6.5
South Africa	5734.633629	Germany	7.0
Zambia	1337.795586	Greece	5.4
Zimbabwe	1445.071062		

264 rows × 1 columns

**Figure 5:** Pandas sections containing the “GDP per capita 2015 (USD)” and “Life Satisfaction Value” for different countries is shown.

The next step of the analysis is to make a final Pandas joint dataframe where are selected only those countries that have available values of both “GDP per capita 2015 (USD)” and “Life Satisfaction Value”. This can be done by using the Pandas join function as shown in details in the accompanying Jupyter Notebook file. At the end we obtain the final Pandas joint data frame as shown in Fig. 6 where only a section of the whole dataframe is shown. The whole joint Pandas dataframe has 38 country entries for the “GDP per capita 2015 (USD)” and “Life Satisfaction Value” values where the country names do not necessarily appear in alphabetic order.

## 4 Data analysis and statistical/machine learning methods

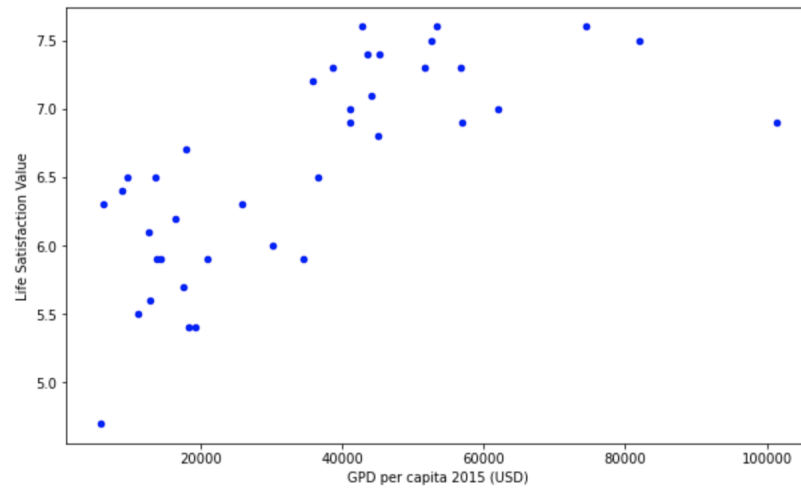
After collecting the data in one dataframe as shown in Fig. 6, the first step to proceed with the analysis is to make a scatter plot and see if there is relation between the data. A scatter plot of the data present in the dataframe of Fig. 6 is shown in Fig. 7 where the feature data ( $X$ ), namely “GDP per capita 2015 (USD)” is plotted versus the output data ( $Y = f(X)$ ), namely “Life Satisfaction value”.

### 4.1 Simple linear regression method (supervised learning)

By looking at the scatter plot in Fig. 7 it seems that the data might follow a linear relationship between them. Indeed, by using the correlation function of the Pandas dataframe, “pd.corr(...)”, a Pearson corre-

	Life Satisfaction Value	GPD per capita 2015 (USD)
Country Name		
Australia	7.3	56755.721712
Austria	7.1	44178.047378
Belgium	6.9	40991.808138
Canada	7.4	43585.511982
Czech Republic	6.7	17829.698322
Denmark	7.6	53254.856370
Finland	7.6	42784.698362
France	6.5	36638.184929
Germany	7.0	41086.729674
Greece	5.4	18167.773727
Hungary	5.6	12706.891215
Iceland	7.5	52564.429179
Ireland	7.0	61995.422803
Italy	6.0	30230.226302
Japan	5.9	34524.469861
Luxembourg	6.9	101376.496574
Mexico	6.5	9616.645006
Netherlands	7.4	45175.231893
New Zealand	7.3	38615.995185
Norway	7.6	74355.515858

**Figure 6:** Pandas section of the joint “GDP per capita 2015 (USD)” and “Life Satisfaction Value” dataframes for different countries is shown.



**Figure 7:** Scatter plot of the “Life Satisfaction Value” vs. “ GDP per capita 2015 (USD)” for 38 different countries of the dataframe in Fig. 6 is shown.

lation coefficient of the value  $r \simeq 0.72$  is found. Clearly such a value of  $r$  indicates that the data might be modelled by a linear relationship of the type  $Y(X) = \beta_0 + \beta_1 X$  where  $\beta_0$  is the intercept value and  $\beta_1$  is the slope value of the linear equation.

In order to find a possible linear relationship, I make use of the “stats” module of the “SciPy” library<sup>2</sup> and use the linear regression built in method:

- `result = sp.stats.linregress(X, Y)`
- `print(result)`  
`LinregressResult(slope=2.399629982572962e-05, intercept=5.741754353755319,`  
`rvalue=0.7202871953226535, pvalue=3.426556470065171e-07, stderr=3.851624914535906e-06,`  
`intercept_stderr=0.15853194959552191)`

The linear regression method used gives the values of the intercept  $\hat{\beta}_0 \simeq 5.74$  and slope  $\hat{\beta}_1 \simeq 2.39 \times 10^{-5}$  with their respective standard errors  $s_{\hat{\beta}_0} \simeq 0.15$  and  $s_{\hat{\beta}_1} \simeq 3.35 \times 10^{-6}$ . The equation for the linear relationship between the data is thus given by the regression line (or least squares line)

$$Y(X) \simeq 5.74 + 2.39 \times 10^{-5} X. \quad (1)$$

In addition to the equation line (1) it is also very useful to have the confidence intervals (CIs) for the regression coefficients  $\beta_{0,1}$  which are give by

$$\beta_{0,1} \in [\hat{\beta}_{0,1} - s_{\hat{\beta}_{0,1}} t_{1-\alpha/2, n-2}, \hat{\beta}_{0,1} + s_{\hat{\beta}_{0,1}} t_{1-\alpha/2, n-2}],$$

where  $t_{1-\alpha/2, n-2}$  is the  $1 - \alpha/2$  percentile (or  $t$ -critical value) of the random  $T$  variable that enters the Student-T distribution function with  $n - 2$  degrees of freedom and  $\alpha$  is the level of significance. If we ask a test statistic for the  $T$  variable at the level of significance of  $\alpha = 0.05$  or confidence level (CL) of 95%, we get for  $n - 2 = 36$  a  $t$ -critical value of  $t_{0.975} \simeq 2.028$ . Thus at the 95% CL, we get the following CIs for the intercept  $\beta_0$  and slope  $\beta_1$  of the OLS linear regression method

$$\beta_0 \in [5.43, 6.04], \quad \beta_1 \in [1.71 \times 10^{-5}, 3.06 \times 10^{-5}].$$

Other two important parameters derived from the linear regression method are the Pearson correlation coefficient,  $r \simeq 0.72$ , that fits the data relatively well and the  $P$ -value of the parameter  $\beta_1$  under the null Hypothesis,  $H_0 : \beta_1 = 0$ . The null Hypothesis is accepted if  $P(\beta_1) \geq 1 - \alpha$  otherwise it is rejected. Since,  $P(\beta_1) \simeq 3.42 \times 10^{-7} \ll 1 - \alpha = 0.95$ , we thus reject the null hypothesis and consequently there is a relation between the data, a fact that is also confirmed from the value of  $r \simeq 0.72$ . In addition we can use the  $t$ -score for the coefficient beta  $t_s = (\hat{\beta}_1 - b)/s_{\hat{\beta}_1}$  to test the null Hypothesis  $H_0 : \beta_1 = b$ . Thus, for  $b = 0$ , we would get a  $t$ -score of

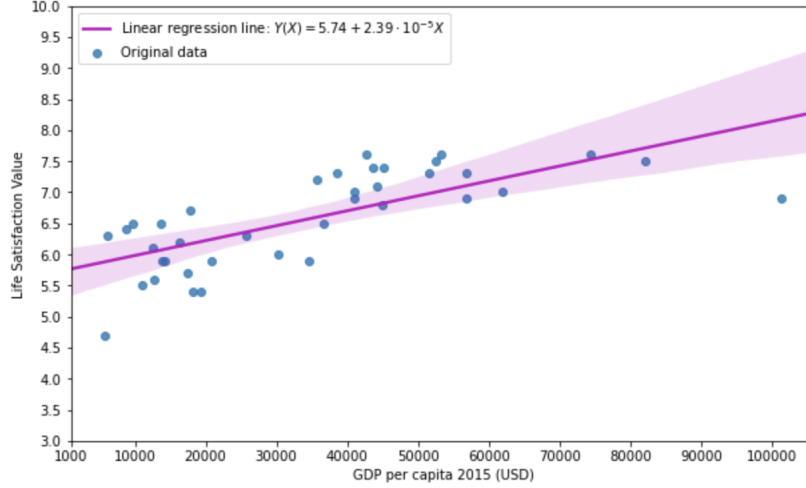
$$t_s = \frac{\hat{\beta}_1 - 0}{s_{\hat{\beta}_1}} = \frac{2.39 \times 10^{-5}}{3.85 \times 10^{-6}} \simeq 6.23.$$

Since  $t_s$  is not in the interval  $[-2.028, 2.028]$ , we reject the null Hypothesis  $H_0 : \beta_1 = b$ . In Fig. 8 the OLS fit line equation (1) and the original data of Fig. 7 are shown.

---

<sup>2</sup>In alternative to the “stats” module one can also use the “sklearn.linear\_model” module of the Scikit-learn library as shown in the accompanying notebook file.





**Figure 8:** Scatter plot of the “Life Satisfaction Value” vs. “ GDP per capita 2015 (USD)” for 38 different countries of the dataframe in Fig. 6 is shown. The region in colour around the regression line represents the 95% confidence band of the least square line. The  $R^2$  value of the data fit is  $R^2 \simeq 0.72$

With the help of the linear relationship in (1) we can make predictions and inference. For example, in the year 2015 it is not known the life satisfaction value of Albania where the country had a GDP per capita of 3952.8 (USD), see the left dataframe in Fig. 5. By using equation (1) and the GDP per capita of Albania,  $X_{\text{Albania}} = 3952.8$  (USD), one can easily find the life satisfaction of Albanians,  $Y_{\text{Albania}} \simeq 5.83$ . On the other hand, we can use equation (1) for inference purposes. For example, it is not known the GDP per capita in 2015 of the countries of Korea and Russia but are known the values of the life satisfaction. Korea in 2015 had a life satisfaction value of  $Y_{\text{Korea}} = 5.9$  while Russia had a life satisfaction value of  $Y_{\text{Russia}} = 5.8$ . By using equation (1), one can find  $X_{\text{Korea}} \simeq 6694.5$  (USD) and  $X_{\text{Russia}} \simeq 2510.4$  (USD).

## 4.2 KNN regression method (supervised learning)

In alternative to the simple linear regression method, one can also use the K-nearest neighbour (KNN) regression method. From the dataset of Fig. 6, we can see that many countries have close values of the “Life Satisfaction Value”. For example, Australia, Canada, Denmark, Finland, Iceland, Netherland, New Zeland, Norway have respectively life satisfaction values of (7.3, 7.4, 7.6, 7.6, 7.5, 7.4, 7.3, 7.6) and the average value of life satisfaction of these countries is 7.4. The KNN regression method can be implemented in python as follows where the default value of regression neighbors is  $K = 5$ .

- `model=sklearn.neighbors.KNeighborsRegressor()`
- `model.fit(X, y)`
- `Xnew=[(input new predictor numerical value)]`
- `print(model.predict(Xnew))`

Now I can make predictions for new life satisfaction values given new GDP per capita values. For



example in our training method, the countries of Albania, United Arab Emirates (UAE) and Armenia were not included in our analysis. These countries had respectively a GDP per capita in the year 2015 (see the whole dataframe on the left in Fig. 5) in USD:  $X_{\text{Albania}} \simeq 3952.8$  (USD),  $X_{\text{UAE}} \simeq 38663.38$  (USD) and  $X_{\text{Armenia}} \simeq 3607.29$  (USD). The KNN regression model predicts the following life satisfaction values:  $Y_{\text{Albania}} = 5.88$ ,  $Y_{\text{UAE}} = 6.98$  and  $Y_{\text{Armenia}} = 5.88$ .

## 5 Data train-test for the linear and KNN regressions with sklearn

In the previous section I analyzed the data by using the whole set of countries present in the dataset in Fig. 6 and made some predictions for the “Satisfaction Value” of countries with no known values of “Satisfaction Value” but with known GDP per capita. However, I did not estimate how accurate are these predictions for future estimates. In this section I perform a data train-test analysis for the linear and KNN regression models in order to estimate how accurate my analysis is in predicting new values of  $Y$  not present in the training dataset. Here, as a matter of example, I split the data of the 38 countries analysed, into a proportion of 80% into the training data and 20% into testing data. I use exclusively the sklearn library for my analysis. For the simple linear regression test-train data, I import the modules and get the values of  $R^2$  score as follows:

- `Xtrain, Xtest, ytrain, ytest=sklearn.model.selection.train_test_split(X, y, test_size=0.2, random_state=0)`
- `model1 = sklearn.linear_model.LinearRegression(fit_intercept=True)`  
`model1.fit(Xtrain, ytrain)`
- `print("Model_1 train  $R^2$  value: ", model1.score(Xtrain, ytrain))`  
`print("Model_1 train  $R^2$  value: ", model1.score(Xtest, ytest))`  
Model\_1 train  $R^2$  value: 0.589117260119626  
Model\_1 test  $R^2$  value: -0.37033956276075175.

In the case of simple linear regression, the train data after splitting performs not as well as previously calculated and gives  $R^2 \simeq 0.589$  in confront of  $R^2 \simeq 0.72$  (whole dataset). On the other hand, the test data linear regression model performs very bad and it gives a negative value of  $R^2 \simeq -0.37$ . This fact indicates that the linear regression method is not much adapted for the test data and it performed worse for the training data.

In the case of KNN-regression model I rescale the data for better numerical stability since the KNN measures the distances of the neighbours (see the accompanying notebook for details) and after I get the values of the test-train  $R^2$  score as follows:

- **from** sklearn.preprocessing **import** StandardScaler  
scaler = StandardScaler()  
 $X_{\text{train\_scaled}} = \text{scaler.fit\_transform}(X_{\text{train}})$   
 $X_{\text{test\_scaled}} = \text{scaler.transform}(X_{\text{test}})$
- model2 = sklearn.neighbors.KNeighborsRegressor(n\_neighbors=3)  
model2.fit( $X_{\text{train\_scaled}}$ ,  $y_{\text{train}}$ )
- **print**("Model.2 train  $R^2$  value: ", model2.score( $X_{\text{train\_scaled}}$ ,  $y_{\text{train}}$ ))  
**print**("Model.2 train  $R^2$  value: ", model2.score( $X_{\text{test\_scaled}}$ ,  $y_{\text{test}}$ ))  
Model.2 train  $R^2$  value: 0.6880950044165277  
Model.2 test  $R^2$  value: 0.7714814814814819

The KNN regression model performs much better than the simple linear regression model and gives better values of  $R^2$  scores, where in the case of the training data,  $R^2 \simeq 0.68$  and in the case of the test data  $R^2 \simeq 0.77$ . So, the KNN model gives a better value of  $R^2$  for the test data and can be used to accurately predict new values not present in the training data. In addition, for the KNN regression model I also calculate the  $MSE$  and  $R^2$  as a function of  $K$  as shown in Fig. 9. As one can see, a value of  $K = 3$  gives the smallest mean square error  $MSE \simeq 0.085$  for the test data and the largest value for  $R^2 \simeq 0.77$ .



**Figure 9:** The mean square error (MSE) and  $R^2$  as a function of the number of nearest neighbor  $K$  is shown. A value of  $K = 3$  gives the best predictions for the test data.

After the analysis done so far, I use the KNN-regression model to predict new values of the “Life Satisfaction Value” for data not present in the test dataset. For example, I consider again the countries of Albania, UAE and Armenia where the GDP per capita is a known quantity. The KNN model after rescaling the GDP per capita data gives  $Y_{\text{Albania}} = 5.8$ ,  $Y_{\text{UAE}} = 7.06$  and  $Y_{\text{Armenia}} = 5.8$ , which is slightly different from what I calculated in the previous section for  $K = 5$  where  $Y_{\text{Albania}} = 5.88$ ,  $Y_{\text{UAE}} = 6.98$  and  $Y_{\text{Armenia}} = 5.88$ . If I had used a value of  $K = 3$  in the previous section, I would get the same values of “Life Satisfaction Value” for the countries of Albania and Armenia and a slightly different value for UAE.

## 6 Take home conclusions

In this article I studied the relationship between the GDP per capita of a given country and the perceived life satisfaction of that country. After analysing the data by using two different regression methods, one is able to predict and make inferences for new data not present in the training set. In se. 4 I evaluated the model performances by calculating the  $R^2$  score for the linear and KNN regression models where the whole dataset of countries present in Fig. 6 was used. No train-test analysis was done. The simple linear regression method, gives a value of the linear correlation coefficient of  $r \simeq 0.72$  and a generalized correlation coefficient of  $R^2 = r^2 \simeq 0.51$ , where  $R^2$  is the generalized correlation coefficient. On the other hand, the KNN regression method gives a generalized correlation coefficient of  $R^2 \simeq 0.69$  for  $K = 5$ . A comparison of  $R^2$  between the two models, would suggest that the KNN regression model would give better fit and predictions for the analysed data.

To asses the model accuracy, in sec. 5, I performed I train-test analysis where the data have been dived in two groups, training data and test data. As shown, the simple linear regression model performed very bad in the test data which suggests that the linear model is not adapted to predict new values. The reason of such low performance is much likely attributed to bias fit and not enough data to make accurate new predictions. On the other hand, the KNN regression model performed quite well in both training and test data, where in the test data a higher value of the  $R^2$  score has been obtained. In addition, I also showed the KNN model performance as a function of  $K$  as shown in Fig. 9, where  $K = 3$  gives the lowest value of  $MSE$  and the highest value of  $R^2$ .