University of Essex

Department of Mathematical Sciences

MA-317: Modelling Experimental Data

# Analysis of the Determinants of Life Expectancy Across the World

Group G:

1900396

1908015

1901094

1900716

1901197

Professors: Dr Joseph Bailey and Dr Stella Hadjiantoni

Date of submission: 11 December 2019

**Abstract**

In this report we investigate the determinants of life expectancy in 2016 using data from the World Development Indicators (WDI). To this aim, a set of linear models are proposed in order to explore the main determinants of life expectancy across countries and an ANOVA analysis to find differences between geographic. We conduct an exploratory statistic analysis and highlight the importance of dealing with the most common issues present in data sets like the one analyzed here. Namely, missing values and collinearity.

*Keywords:* life expectancy, missing values, collinearity, linear regression, ANOVA

**Contents**

**Word Count**: 3,366

# 1. Introduction

In this report, the determinants of life expectancy are explored taking many country-level aggregated factors into consideration. The relevance of this analysis is that understanding can give an insight into many other well-being and degree of development indicators. In this way, public policies impact can be assessed.

To this extent, using World Bank's *World Development Indicators* data, we propose a set of linear models that can be used to explore the main determinants of life expectancy across countries and, at the same time, use them to predict next years' tendencies.

Along the report, we conduct exploratory statistic analysis and highlight the importance of dealing with the most common issues present in data sets like the one analyzed here. The rest of the report is organized as follows: in section 2, a *descriptive statistics analysis* of the data is conducted, giving special attention to the relation with the variable of interest. Next, several procedures to deal with *missing data* are deployed. Then, we test for *collinearity* between explanatory variables. In section 3.1, four different linear models are proposed to identify the determinants of the life expectancy and one-way ANOVA analysis across continents is conducted. Finally, section 4 concludes.

# 2. Preliminary Analysis

In this section we describe the data and provide descriptive statistics. Additionally, we propose three different methods to deal with missing data based on the distribution characteristics of each variable, and highlight the disadvantages of conducting a complete case analysis. Finally, we test for collinearity between some variables based on different categories.

## 2.1. Dataset and Descriptive Statistics

The dataset was obtained from the World Bank's *World Development Indicators* (WDI) compilation. This data comes from officially-recognized international sources. The dataset contains 22 variables, being life expectancy the variable of interest, with a total number of 217 observations corresponding to recognized countries for 2016. Additionally, there are 47 observations of regional and economic aggregations such as; European Union and OECD members.

Table 1 shows a description of all the variables. For life expectancy, there are 199 observations (i.e. 18 missing values), with a minimum of 51 years, a maximum of 84 years and a average life expectancy of 72.3 years with a standard deviation of 7 years. A definition for each variable can be seen in appendix A.

We can observe that there are some variables (for instance, unsafe water mortality, population growth, secondary education, health expenditure per capita, GDP per capita, and mobile subscriptions) have very high dispersion. This can represent an issue, as their scale and variability may affect

the correlation magnitude with life expectancy. Therefore, they are log transformed for the rest of the analysis to make the relationship between the variables and life expectancy close to being symmetric. Also, it would help in straightening out the data and improve the regression model.

**Table 1:** Descriptive statistics

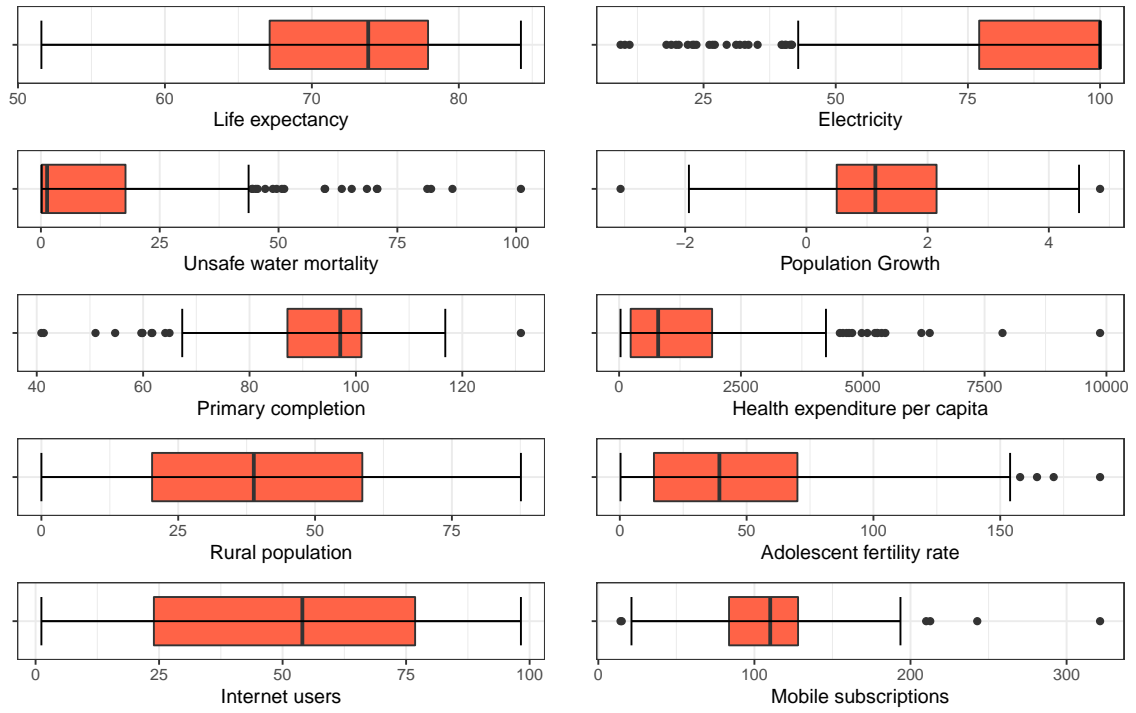|  | Mean | Median | Standard Deviation | Min. | Max. | NA.s |
|---|---|---|---|---|---|---|
| Life Expectancy | 72.30 | 73.84 | 7.78 | 51.59 | 84.23 | 18 |
| Electricity | 84.86 | 99.99 | 25.42 | 9.30 | 100.00 | 2 |
| Adjusted Income | 355.73 | 29.78 | 1450.69 | 0.16 | 15985.10 | 45 |
| Children out of school | 6.27 | 1.82 | 9.21 | 0.00 | 42.62 | 85 |
| Primary education exp. | 31.31 | 30.57 | 10.66 | 12.70 | 64.06 | 146 |
| PPP | 341.76 | 211.35 | 331.19 | 88.90 | 950.17 | 211 |
| Unsafe water mortality | 12.50 | 1.30 | 20.82 | 0.10 | 101.00 | 34 |
| Adult literacy rate | 90.21 | 94.65 | 14.70 | 22.31 | 99.99 | 183 |
| Population growth | 1.29 | 1.14 | 1.22 | -3.07 | 4.85 | 1 |
| Population total | 34.26 | 6.42 | 134.80 | 0.01 | 1378.66 | 1 |
| Primary completion | 92.10 | 97.06 | 15.71 | 40.87 | 131.02 | 86 |
| Secondary ed. | 6.37 | 6.00 | 0.92 | 4.00 | 9.00 | 13 |
| Secondary ed. teachers | 209.94 | 42.12 | 721.75 | 0.04 | 6219.58 | 94 |
| Health exp. | 6.74 | 6.27 | 2.99 | 1.75 | 23.29 | 31 |
| Health exp. per capita | 1426.02 | 801.76 | 1694.91 | 29.91 | 9869.74 | 33 |
| Unemployment | 8.34 | 6.29 | 6.21 | 0.15 | 27.47 | 106 |
| Youth Unemployment | 19.39 | 15.96 | 12.58 | 0.49 | 54.31 | 115 |
| Rural population | 39.68 | 38.81 | 24.07 | 0.00 | 87.61 | 3 |
| Adolescent fertility rate | 48.09 | 39.26 | 40.57 | 0.29 | 189.38 | 23 |
| GDP per capita | 20698.04 | 13247.65 | 21828.44 | 743.90 | 123573.63 | 24 |
| Mobile subscriptions | 107.21 | 110.14 | 40.37 | 14.25 | 321.45 | 16 |
| Internet users | 51.32 | 54.00 | 28.96 | 1.18 | 98.24 | 13 |

For the rest of the analysis we focus on only those variables that are highly correlated with life expectancy (see figure *scatter*). In this way, we choose the ten variables with the highest correlation coefficients, their corresponding boxplots are displayed in figure *boxplot*. This was so, as to maintain a simpler analysis. These variables are: electricity, unsafe water mortality, population growth, primary completion, health expenditure per capita, rural population, adolescent fertility rate, internet users, and mobile subscriptions.

We can see in figure 1 that some distributions are very skewed, e.g. electricity, unsafe water mortality and health expenditure per capita. Additionally, some have very high variance such as; internet users.

The relationship between life expectancy and GDP per capita in logarithm form was however analyzed with more detail, as we would expect it to be one of the principal determinants of life expectancy. Consequently, from the scatter plot in figure B.3 (see appendix B) we will expect that *log(GDP pc)* should explain, at least, 67.9% of variance of life expectancy.

### 2.2. Handling Missing Data

In resolving missing values which mostly occur in statistical analysis, different methods can be used. Amongst these include the *complete case analysis*. Here, observations with missing values for the explanatory variables are omitted, leaving only complete cases. Nevertheless, this could lead to

**Figure 1:** Distribution of variables of interest

biased estimators if the data is not missing completely at random (MCAR). However, even when this alternative is likely to be unbiased, it discards most of the information contained in the variability of the data. Thereby, is the least effective method in resolving missing values [1]. In this regard, keeping only strictly complete cases drops the number of observations to 0. Since there were no countries with complete information. For this reason, the conditions for complete cases were relaxed.
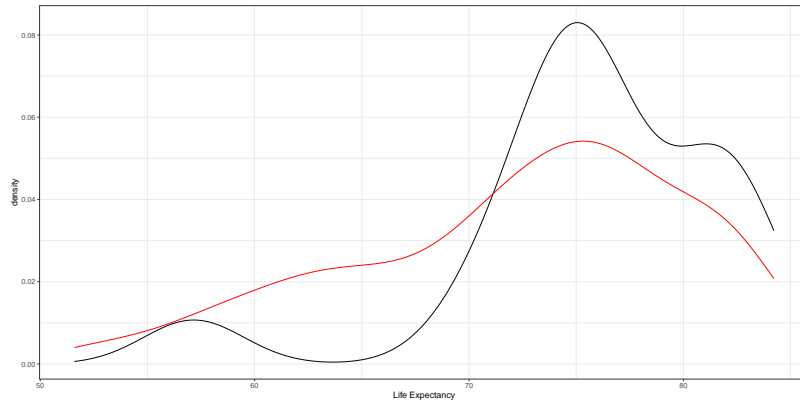
Out of a total of 217 countries, PPP, Literacy rate and Primary expenditure had 211, 183 and 146 missing values respectively (table 1). Thus, they were removed to relax the conditions. This then left us with 48 observations. We then assumed this as a complete case (i.e. excluding PPP, literacy rate and Primary education expenditure).

In the figure 2, the probability density function (PDF) of life expectancy under the *assumed complete case* was then plotted alongside with the PDF of life expectancy considering all observations (black and red curves, respectively).

From the density functions we can infer that countries with lower life expectancy are underrepresented as the distribution is skewed towards the left (i.e. its mass is concentrated in countries with higher life expectancy). Moreover, if we keep only those with complete observations, many of this observations are lost, making our analysis biased as the missingness appears to be not at random.

As we shown in figure B.3, life expectancy is closely correlated with the level of economic development, so we could expect that in lower income countries statistical data monitoring is less efficient.

To avoid the bias that would result from a complete case analysis, different predictor variables and

**Figure 2:** Probability Density Function for Life Expectancy

their features gathered from table 1 gave an insight on what method of imputation is the more optimal.

The imputation methods deployed include: i) Mean Imputation of missing values; ii) Median Imputation of missing values; iii) Linear Regression Method to perform deterministic imputation

Section 2.1 showed that the predictor variables and number of missing values. The criteria used to select the imputation method was the number of missing values. Therefore, *mean imputation* was used to deal with those variables with less than 100 missing values and symmetric distributions (i.e. similar median and mean values). In contrast, in the case of variables with less than 100 missing values but skewed distributions (i.e. large difference between mean and median values), *median imputation* was performed, as using the mean would drastically change their distribution.

Lastly, from the table 1, values within the range of 100-200 missing values were imputed with the *linear regression method* for deterministic imputation. As there are many non-observed data points, implementing the two previous methods would result in a very concentrated distribution around the mean (i.e. very high kurtosis), then most of the information provided by these variables would be lost.

*Public Private Partnership Investment* had a total of 211 NA's (that is, missing data) out of 217 countries. This is rather insignificant as there are too many missing values to help determine which method of imputation is to be deployed. Hence, inference would be not statistically reliable, therefore this variable was omitted for the rest of the analysis.

After imputation, the PDF of all imputed predictor variables were plotted to see the difference before and after imputation was made (see figure B.2 in appendix B).The red and black functions represent the distribution before and after imputation, respectively. Thus, in general, there were no major changes in their distributions.

### 2.3. Collinearity

In regression analysis, it can be the case that an explanatory variable could be expressed as a linear combination of the rest. In this case, we can say we are in the presence of collinearity. This issue

signifies a problem, analytically the matrix $X^t X$, becomes singular (i.e. it is not invertible) and the parameters of the model can not be derived. In practice, they significantly increase the variance of the estimators. For this reasons, we conducted an analysis in order to determine if there is collinearity in our dataset.

There are several methods to establish collinearity. One of those consists in conducting a correlation test, to determine if the relation between two given variables is statistically significant. Another method consists in estimating the Variance Inflation Factor (VIF), which is computed from the $R_j^2$ resulting from regressing variable $X_j$ against the reamining ones, or a subset that we suspect may be correlated with it. VIF can be defined as follows:

$$VIF = \frac{1}{1 - R_j^2} \tag{2.1}$$

There is no analytical rule on which is the correct threshold for this factor. However, some authors propose a maximum value 5 [2, 3]. In the present analysis, we set our threshold to this value. That is, if VIF$> 5$ we can say that the variables involved are linearly correlated.

**Table 2:** Results of Collinearity Analysis

| | *Dependent variable:* | | | | | |
|---|---|---|---|---|---|---|
| | log(GDP per capita) | | | | | log(Water Mort.) |
| | (1) | (2) | (3) | (4) | (5) | |
| Rural Population | $-0.006^{***}$ | $-0.032^{***}$ | | | | |
| | (0.001) | (0.002) | | | | |
| log(Health Exp. pc) | $0.684^{***}$ | | $0.807^{***}$ | | | $0.009^{***}$ |
| | (0.049) | | (0.029) | | | (0.001) |
| Electricity | $0.003^{**}$ | | | $0.032^{***}$ | | |
| | (0.002) | | | (0.002) | | |
| Internet Users | $0.003^{*}$ | | | | $0.033^{***}$ | |
| | (0.002) | | | | (0.002) | |
| Health Exp. | $-0.103^{***}$ | | | | | |
| | (0.011) | | | | | |
| Adjusted Income | 0.009 | | | | | |
| | (0.016) | | | | | |
| Unemployment | $-0.011^{**}$ | | | | | |
| | (0.005) | | | | | |
| log(Mobile Subs.) | 0.115 | | | | | |
| | (0.078) | | | | | |
| Constant | $4.702^{***}$ | $10.647^{***}$ | $4.088^{***}$ | $6.619^{***}$ | $7.677^{***}$ | $4.129^{***}$ |
| | (0.494) | (0.107) | (0.191) | (0.181) | (0.088) | (0.051) |
| Observations | 217 | 217 | 217 | 217 | 217 | 217 |
| Correlation test | 0.654 | -0.690 | 0.887 | 0.735 | 0.830 | -0.834 |
| $R^2$ | 0.889 | 0.476 | 0.787 | 0.540 | 0.688 | 0.334 |
| VIF | 9.008 | 1.909 | 4.688 | 2.173 | 3.206 | 1.502 |

*Note:* $^{*}$p$<$0.1; $^{**}$p$<$0.05; $^{***}$p$<$0.01

Collinearity is tested within two categories: economic and health. We chose this categories as we would expect the variables in each to be correlated either because they are measuring the same thing

or there is a third confounder non-observed establishing a mechanism of relation between them.

As seen in table 2, rural population, electricity and internet users have VIF values of 1.90, 2.17 and 3.21, respectively, on GDP pc. Therefore, the relation is not strong enough to establish collinearity. In the same way, the value of VIF is roughly 1.50 in the health category, which is much lower than 5.

In contrast, for health expenditure per capita, the value of VIF is approximately 4.69 (very close to the threshold) which imply a potential collinearity problem with GDP pc (this relation can be explained by the government expenditure multiplier). Finally, model 1 shows that GDP pc appears to be highly correlated with a set of selected variables, with a VIF equal 9.01. Hence, it can be said we are in the presence of collinearity.

As this issue increases the variance of the linear model estimators (i.e. high variance), there are three widely used methods to deal with it: i) Ignore collinearity if it is not strong enough; ii) Use an alternative estimation method such as ridge regression [4]; iii) Implementing the variable selection technique to discard the correlated variable.

## 3. Analysis

Life expectancy at birth can be explained by several factors specific to each country. Namely, quality and coverage of health services, other public services like security, infrastructure and education, eating patterns, genetics, and even geographical factors (e.g. local outbreaks, risk of natural disasters, etc.) [5]. In this way, a great number of studies establish a correlation between life expectancy and a set of economic indicators that can help to explain its evolution through time and therefore, the source of variations across countries [5, 6, 7, 8].

In this section, we try to identify the determinants of life expectancy at an aggregated country level. The results will help to understand the role of economic development and how public policies affect life quality. To achieve this, a linear regression and ANOVA analyses are conducted.

### 3.1. Regression Analysis

We propose a multivariate linear model to explore which variables in our data set are good predictors of life expectancy. One of the advantages of using this type of models is the parameters are easy to interpret and, at the same time, it can be used to model non-linear relations between explanatory variables and the dependent variables by transforming the data, maintaining the linearity of the parameters (e.g. log, polynomial or exponential-transformations).

From the analysis in section 2.1 (see figure B.1 in the appendix) we decide to initially consider only those variables with the highest correlation with life expectancy in order to keep the analysis parsimonious. The criteria used to determine relevance is the *Pearson correlation coefficient* to be greater than $|0.5|$ (the correlation matrix is shown in figure B.5 in the appendix).

In section 2.3 it was found that *log(GDP pc)* is highly correlated with other variables in the dataset. The reason why we tested collinearity between these is because we would typically use those variables to explain the GDP per capita [9]. In this context, including the variable *GDP per capita* would be redundant to the analysis, because it adds no additional explanatory power to the model and, in fact, would increase the variance of the estimators as shown by the VIF in table 2.

For this reason, two different saturated models are proposed. The first one excludes *GDP pc*, and the second one includes this variable and a series of controls.

**Saturated Model**

$$life\ expectancy = X\beta + \epsilon \tag{3.1}$$

where *life expectancy* $\sim N(X\beta, \sigma^2)$ is the dependent variable, $\epsilon \sim N(0, \sigma^2)$ a random error term, $\beta$ the vector of parameters to be estimated and $X$ a vector of covariates that includes: *electricity*, *log(water mortality)*, *population growth*, *log(adjusted income)*, *primary completion*, *health expenditure pc*, *health expenditure pc$^2$*, *rural population*, *adolescent fertility rate*, *internet users* and *mobile subscriptions*.[1]

As the relation between *life expectancy* and *health expenditure per capita* appears to be non-linear we decide to include a polynomial term of second degree for this variable (see figure B.1 in the appendix). The relation between these two variables is positive at a diminishing rate (concave shape), for this reason we would expect the linear term to be positive and the squared term to be negative.

**Log-log model**

$$log(life\ expectancy) = \beta_0 + \beta_1 log(GDP\ pc) + \beta_2 Gini\ Index + \beta_j contintent_j + \epsilon \tag{3.2}$$

The dependent variable and error term have the same assumptions as in the previous model, in this case the explanatory variables are: *log(GDP pc)*, *Gini index* and a continent indicator (Africa is the omitted group[2]). The parameters of a *log-log* model can be interpreted as the elasticity of the dependent against the explanatory variable.[3]

*3.1.1. Results*

In order to select the variables that lead to the smaller expected squared error, we perform *stepwise* variable selection. As there is no hard rule on which information criterion is the best, both the *Akaike Information Criterion* (AIC) and *Bayesian Information Criterion* were considered.

After using the *step* command in R in both directions, the results are shown in table 3. Model 1 shows the best model starting from the saturated model. All of the variables are statistically significant and *log(water mortality)* is the one with the highest effect on life expectancy. Nonetheless, as life

---

[1]A description of each variable is given in appendix A.

[2]Every continent parameter is interpreted as the difference with respect to the omitted group.

[3]For small changes, shows the percentage change of the dependent variable when the explanatory changes by 1%.

expectancy is estimated from the mortality rates at birth of the corresponding cohort, as if they would keep constant throughout the life of the cohort individuals, including the water mortality rate is trivial [10].

For this reason, we omit this variable and perform *stepwise* variable selection again. Model 2 is the one selected based on the AIC and model 3 on the BIC.

In model 2 we observe that the principal determinant of life expectancy is the health expenditure per capita. In this way, life expectancy could be used as an indicator of the public health policy efficacy, however higher mortality would in fact lead to an increase in health expenditure. Other significant explanatory variables are electricity coverage and adolescent fertility. According to the World Health Organization, complications during pregnancy and childbirth are the leading cause of death for teenage girls [11]. Additionally, the access to electricity services and life expectancy correlation could be explained by a hidden variable that reflects the living standards quality.

Model 3 reflects the fact that BIC penalizes more heavily in the number of parameters to estimate, in this way the best model is the one without the squared term.

Finally, in model 4 all the parameters are statistically significant. As mentioned before, the *GDP pc* parameter can be interpreted as an elasticity, in this way a 1% increase in the mean GDP per capita would lead to an increase of 0.05% in the mean life expectancy. But, it is important to note that this correlation does not imply causality, since the relation may be endogenous (i.e. they are mutually correlated). Regarding the continent indicators, every continent has a positive correlation with life expectancy compared to Africa (the omitted group) (see figure B.6).

In order to assess linear regression assumptions fulfillment, a discussion is presented in the appendix C.

So as to determine what specification is the best to predict new unseen observations, cross-validation techniques could be implemented (i.e. randomly splitting the data into test and training sets). Another solution is testing our models with different years observations, or by predicting the life expectancy missing values in our data set and validate these predictions with a secondary source.

### 3.2. ANOVA

A One-Way ANOVA (Analysis of Variance) is a statistical technique by which we can test if three or more means are equal, i.e. if one differs significantly among three or more levels of a factor.

The benefits of using this method are: it is easy to estimate and can be manually computed using simple algebra rather than complex matrix calculations. It can control the overall Type I error rate and provide an overall test of equality of group means. If normality assumption is true then the test is more powerful.

**Table 3:** Results of Regression Analyses

| | Dependent variable: | | | |
|---|---|---|---|---|
| | Life Expectancy | | | log(Life Expectancy) |
| | (1) | (2) | (3) | (4) |
| Electricity | 0.101*** | 0.120*** | 0.110*** | |
| | (0.017) | (0.016) | (0.015) | |
| log(Water Mort.) | −0.835*** | | | |
| | (0.259) | | | |
| Health Exp. pc | 0.629*** | 1.751*** | 0.705*** | |
| | (0.209) | (0.556) | (0.213) | |
| Health Exp. pc$^2$ | | −0.136** | | |
| | | (0.067) | | |
| Adolscent Fertility | −0.028*** | −0.036*** | −0.041*** | |
| | (0.010) | (0.009) | (0.009) | |
| Internet Users | 0.060*** | 0.074*** | 0.084*** | |
| | (0.018) | (0.018) | (0.016) | |
| Mobile Coverage | −1.179* | −1.049 | | |
| | (0.690) | (0.699) | | |
| log(GDP pc) | | | | 0.053*** |
| | | | | (0.005) |
| Gini Index | | | | −0.003*** |
| | | | | (0.001) |
| Asia | | | | 0.084*** |
| | | | | (0.013) |
| Europe | | | | 0.103*** |
| | | | | (0.016) |
| North America | | | | 0.145*** |
| | | | | (0.019) |
| Oceania | | | | 0.091*** |
| | | | | (0.018) |
| South America | | | | 0.105*** |
| | | | | (0.014) |
| Constant | 67.143*** | 63.250*** | 59.834*** | 3.811*** |
| | (3.048) | (2.859) | (1.464) | (0.043) |
| Observations | 195 | 195 | 195 | 195 |
| $R^2$ | 0.834 | 0.828 | 0.823 | 0.772 |
| AIC | **460.81** | **467.05** | 469.37 | -1,126.7 |
| BIC | 484.34 | 490.58 | **486.18** | -1,100.5 |

*Note*: *p<0.1; **p<0.05; ***p<0.01

In this section ANOVA analysis is conducted to test if there is a statistical difference of mean life expectancy between the 5 different continents. The M49 classification from the United Nations is used.

The results are shown in table 4, the five continents average of life expectancy differed significantly on anxiety level, $F_{(5,189)} = 2e - 16$,p<0.05. As a consequence, this difference can be a reflect of economic and political trends across regions. The p value on the table is 2e-16 which is less than 0.05 indicates that two or more groups have significantly different means.

### 3.3. Define Null and Alternative Hypotheses

Life expectancy of $\mu_{Asia} = \mu_{Oceania} = \mu_{NorthAmerica} = \mu_{SouthAmerica} = \mu_{Europe}$

Not all of $\mu$ are equal. Alpha: $\alpha = 0.05$
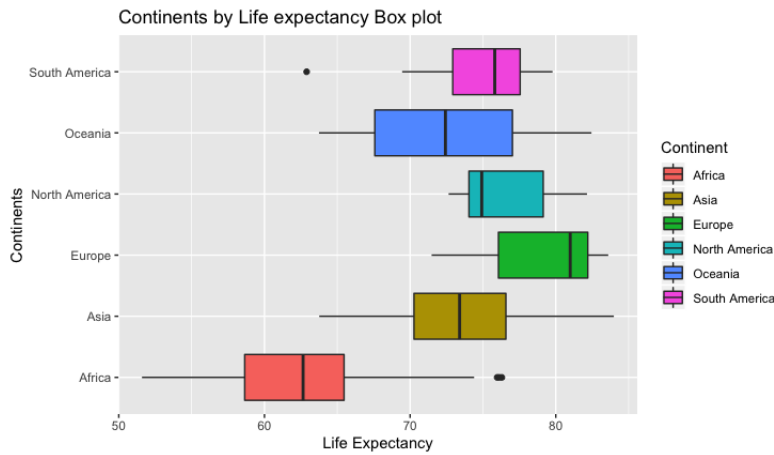
### 3.4. Results

**Table 4:** ANOVA test

|            | Df  | Sum sq | Mean sq | F value | $Pr > F$     |
|------------|-----|--------|---------|---------|--------------|
| Continents | 5   | 7016   | 1403.2  | 57.92   | $< 2e - 16$  |
| Residuals  | 189 | 4579   | 24.2    |         |              |

**Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05'.' 0.1 ' ' 1**
**18 observations deleted due to missingness**

If F is greater than $F_{(0.05,5,189)} = 2.261892$ reject the null hypothesis.

The five continents average of life expectancy differed significantly on anxiety level, $F_{(5,189)} = 2e - 16$,p<0.05.

The p value on the table is 2e-16 which is less than 0.05 indicates that two or more groups have significantly different means.



**Figure 3:** This box chart shows that there is notable different in each sample. Africa's life expectancy is lower than others continents and have a huge differ life expectancy between max and min. The life expectancy in Asia and Oceania are near normal distribution. The box plot of North America is left-skewed. In contrast, Europe and South Americas' are right-skewed.

A discussion on the assumptions of ANOVA is presented in appendix D.

## 4. Discussion

Understanding the determinants of life expectancy could be useful to explore the extent to which its difference across countries and continents is explained by environmental random factors and by economic and political ones. This in turn, could be used by international organizations to assess the effect and impact of such policies. Namely, if public expenditure is being allocated efficiently, how important is education to improve people's well-being, among others. In this way, helping into the design of better plans and social programs.

However, this type of data is prone to have several statistical issues. For instance, the quality of the information collection techniques varies substantially across countries. In fact, those with lower

a lower level of development tend to have incomplete data. As a consequence, it is important to correctly identify, evaluate and correctly deal with problems like: missing data, collinearity, unbalanced categories, high variance variables, etc.

Throughout this report a statistical preliminary analysis was conducted, including: descriptive sattistics to identify the properties of the data and possible issues, methods to deal with missing data and collinearity test. Besides, a series of lineal models were proposed aimed to better predict and understand life expectancy, as well as an ANOVA analysis to test difference across continents.

It was found that the most important factors are those related to health services and sanitary factors, and the degree of access to basic services as electricity services and education. Moreover, there is a significant difference across regions. It is important to note that these relations does not imply causality, as it may be endogenous. To this end, further analysis through identification methods like instrumental variables could be explored.

**Contributions**: All members contributed to the elaboration of the report. However, each one focused on a specific task:

- **1908015**: Task 1
- **1901197**: Task 2
- **1900716**: Task 3
- **1900396**: Task 4
- **1901094**: Task 5

## References

[1] Sterne AC, White I, Carlin J, Spratt M, Royston P, Kenward M, and Carpenter A, Woodand J. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *Bmj*, 338: b2393, 2009.

[2] Rogerson PA. Data reduction: factor analysis and cluster analysis, 2001.

[3] Shah MH and Samdani S. Impact of trade liberalization on fdi inflows to d-8 countries. *Shah, MH, & Samdani, S.(2015). Impact of trade liberalization on FDI inflows to D-8 countries. Global Management Journal for Academic & Corporate Studies (GMJACS)*, 5(1):30–37, 2015.

[4] Draper NR and Smith H. *Applied regression analysis*, volume 326. John Wiley & Sons, 1998.

[5] Rodgers G. B. Income and inequality as determinants of mortality: An international cross-section analysis. *Population Studies*, 33(2):343–351, 1979. ISSN 00324728.

[6] Wilkinson R. G. Income distribution and life expectancy. *BMJ*, 304(6820):165–168, 1992. ISSN 0959-8138.

[7] Messias E. Income inequality, illiteracy rate, and life expectancy in brazil. *American journal of public health*, 93:1294–6, 09 2003.

[8] Shaw J., William H., and Ronald V. The determinants of life expectancy: An analysis of the oecd health data. *Southern Economic Journal*, 71(4):768–783, 2005. ISSN 00384038.

[9] Parkin D., McGuire A., and Yule B. Aggregate health care expenditures and national income: is health care a luxury good? *Journal of health economics*, 6(2):109–127, 1987.

[10] OECD. *Health at a Glance 2019*. 2019.

[11] World Health Organization. Global health estimates 2015: deaths by cause, age, sex, by country and by region, 2000–2015, 2016. Geneva.

**Appendix**

## A. Glossary

**Life Expectancy**: Life expectancy at birth, total (years)

**Electricity**: Access to electricity (% of population)

**Adjusted income**: Adjusted net national income (current US$)

**Children out of school**: Children out of school (% of primary school age)

**Primary education expenditure**: Expenditure on primary education (% of government expenditure on education)

**PPP**: Public private partnerships investment in water and sanitation (current US$)

**Unsafe water mortality**: Mortality rate attributed to unsafe water, unsafe sanitation and lack of hygiene (per 100,000 population)

**Adult literacy rate**: Literacy rate, adult total (% of people ages 15 and above)

**Population Growth**: Population growth (annual %)

**Population total**: Population, total

**Primary completion**: Primary completion rate, total (% of relevant age group)

**Secondary ed.**: Secondary education, duration (years)

**Secondary ed. Teacher**: Secondary education, teachers

**Health expenditure**: Current health expenditure (% of GDP)

**Health expenditure per capita**: Current health expenditure per capita, PPP (current international $)

**Unemployment**: Unemployment, total (% of total labor force) (national estimate)

**Youth unemployment**: Unemployment, youth total (% of total labor force ages 15-24) (national estimate)

**Rural population**: Rural population (% of total population)

**Adolescent fertility rate**: Adolescent fertility rate (births per 1,000 women ages 15-19)

**GDP per capita**: GDP per capita, PPP (current international $)

**Mobile subscriptions**: Mobile cellular subscriptions (per 100 people)

**Internet users**: Individuals using the Internet (% of population)
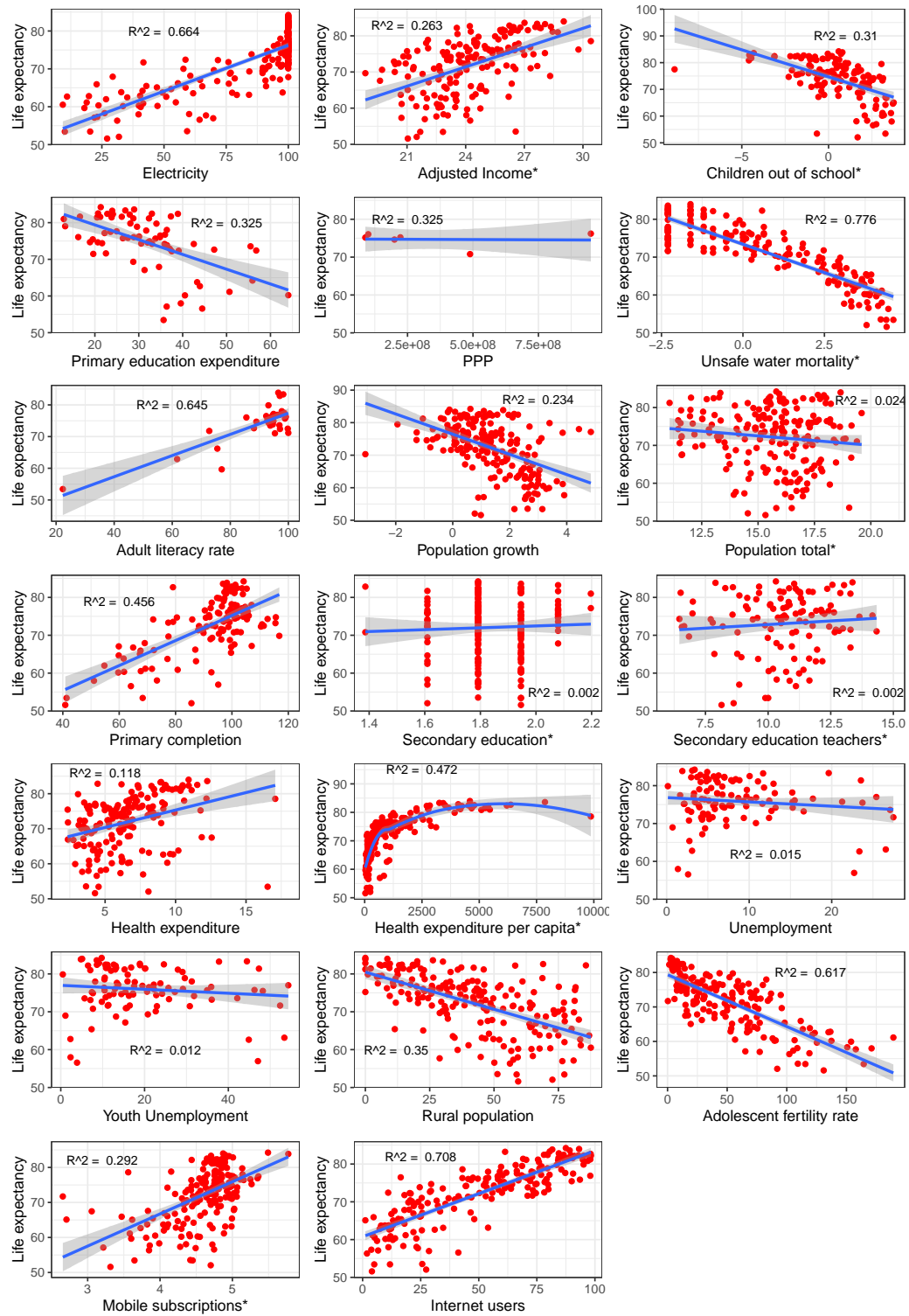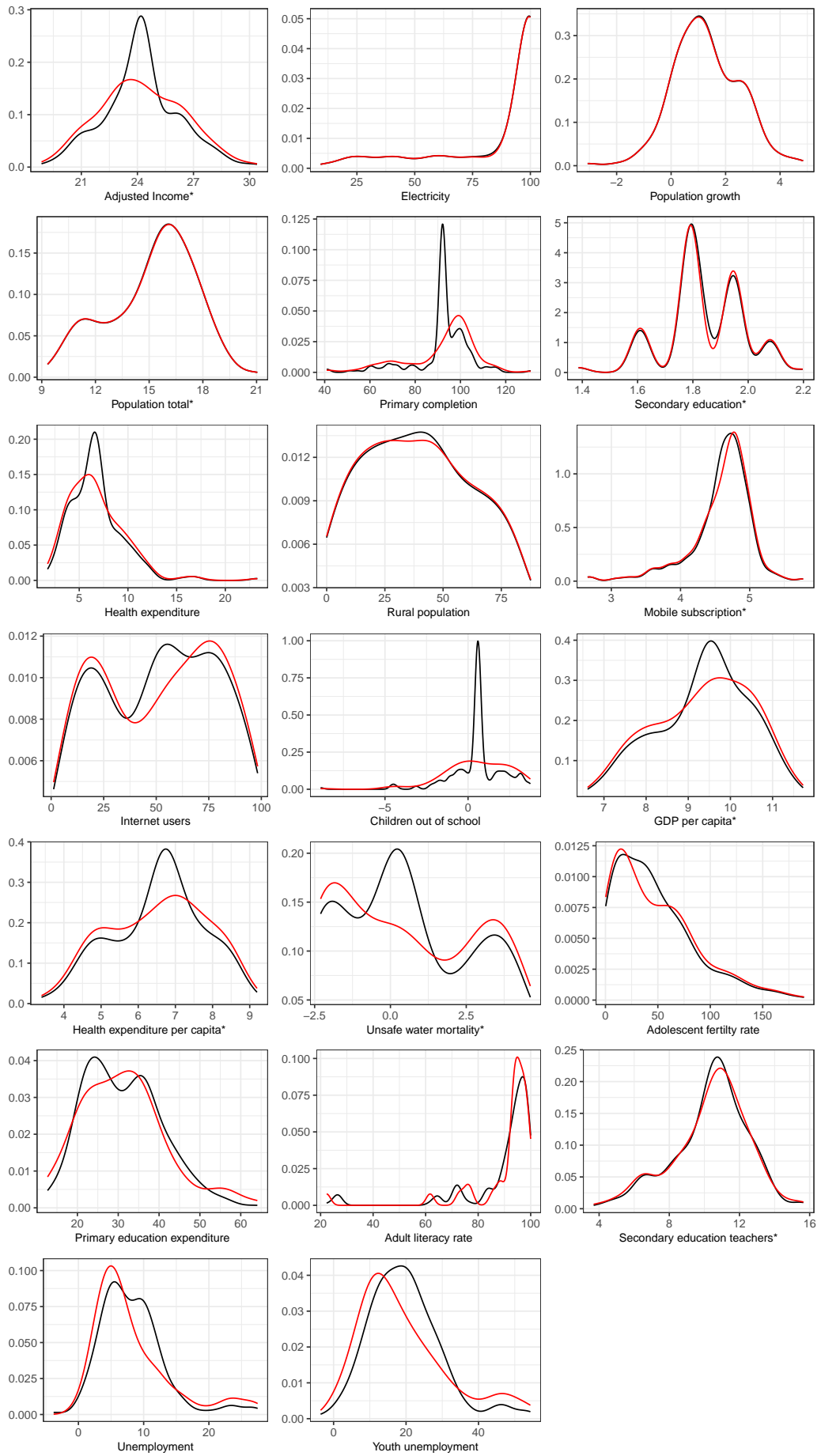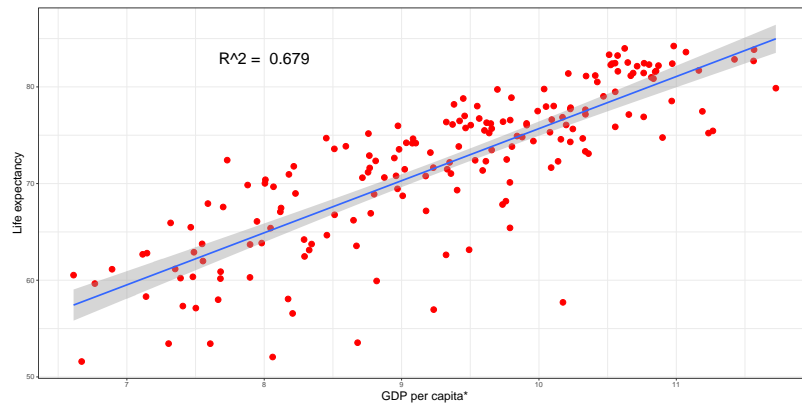
## B. Additional Plots



**Figure B.1:** Life Expectancy against all the Variables
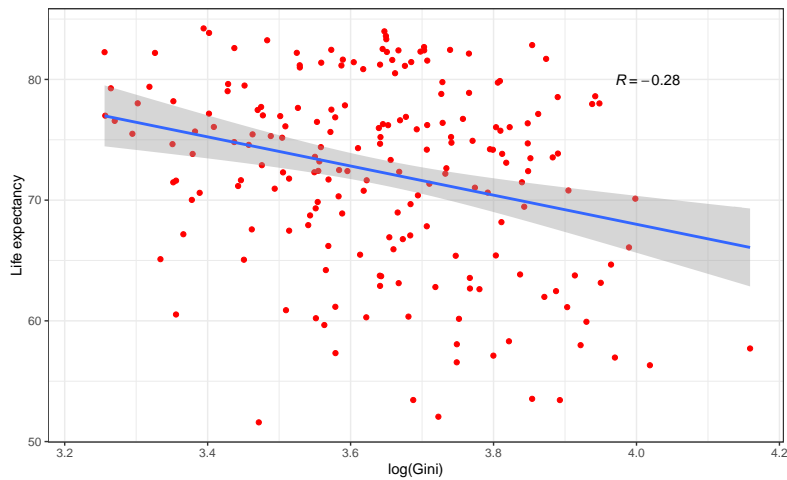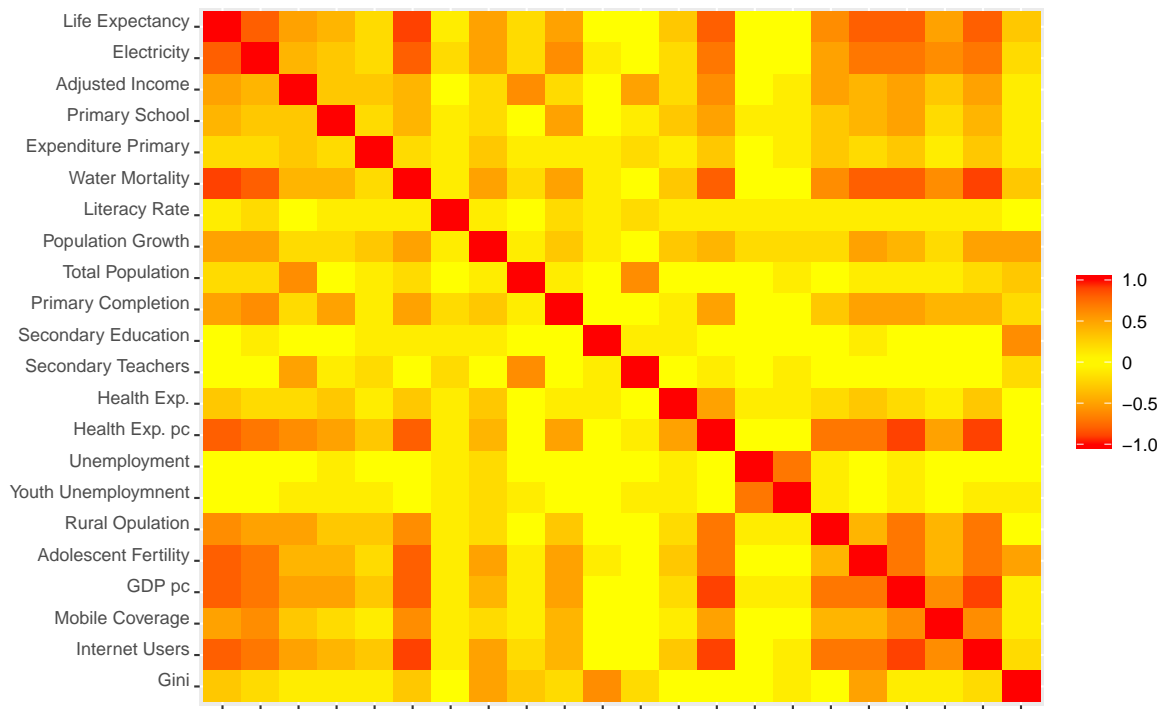
**Figure B.2:** Probability Density Function for Variables
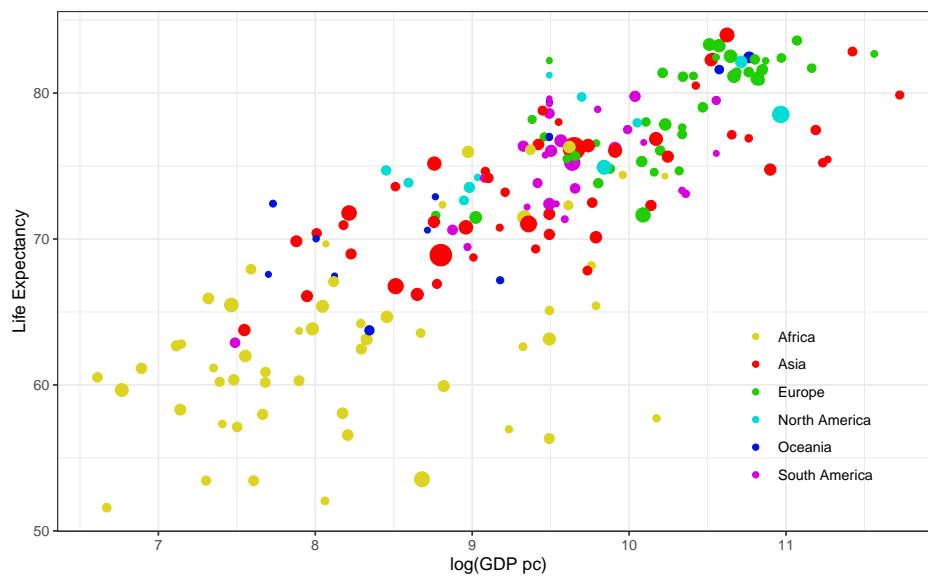
15

**Figure B.3:** Correlation between Life Expectancy and log(GDP pc)



**Figure B.4:** Gini Index vs Life Expectancy



**Figure B.5:** Correlation Matrix for the Dataset

**Figure B.6:** Life Expectancy by Contnitnet

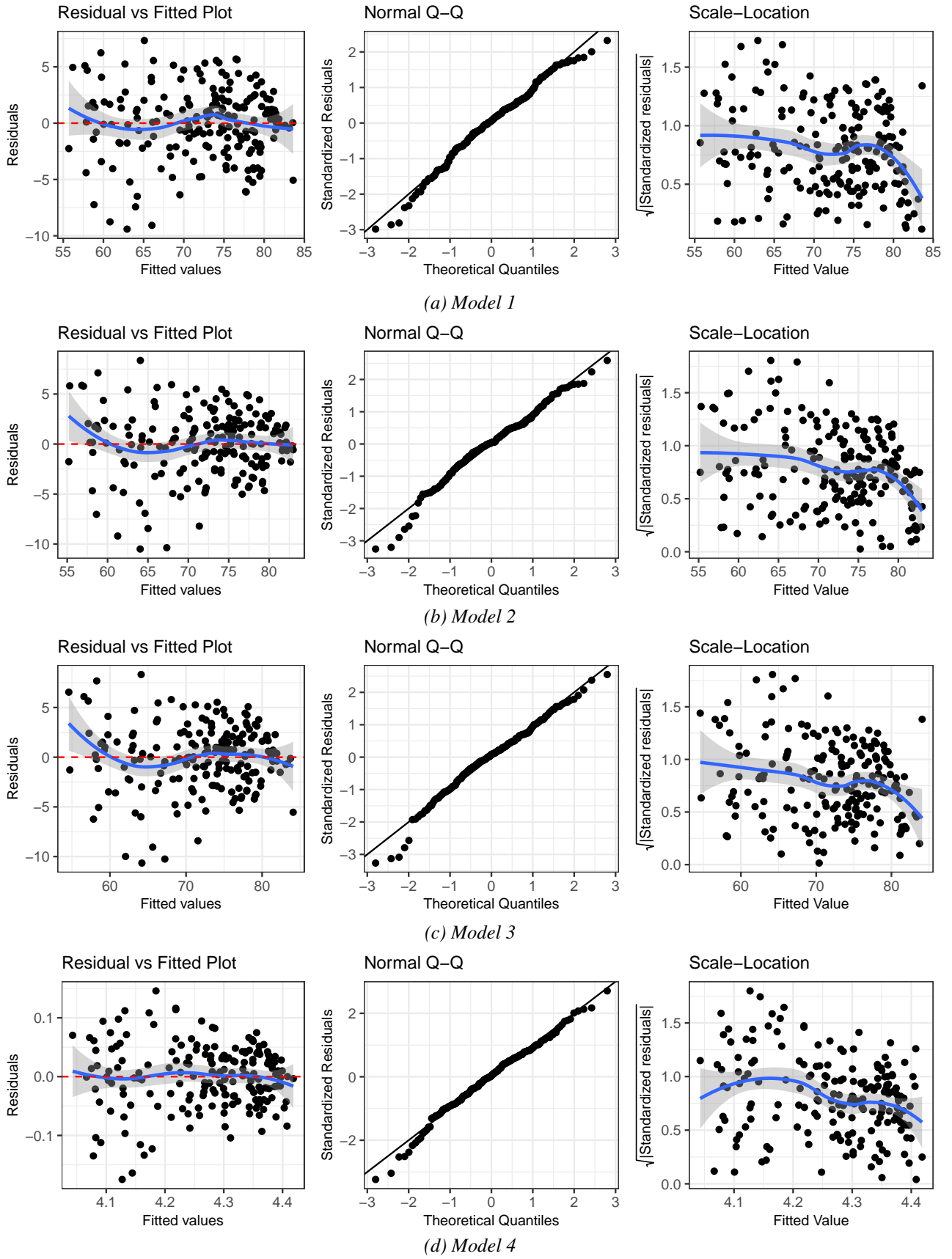## C. Discussing the Assumptions of Linear Regression Analysis

The *residual vs fitted plot* shows the linearity assumption between dependent and explanatory variables, we would expect the errors to be randomly and equally distributed around zero across all the range of life expectancy. In all the models this assumption appears to hold for those values around mean life expectancy. However, the lower age range seem to be underrepresented in the sample. Another reason may be that the determinants of life expectancy are different along the distribution. To solve this issue a quantile regression could be performed, in this way we would be able to find how the weights of each explanatory variable change across the distribution.

The *Q-Q plot* plot is used to assess visually the normality assumption. It appears to fulfill in general, but, as in the previous paragraph, distribution tails are not well-behaved. The reason could be that, as life expectancy and national income are positively related, lower income countries may have poor quality data.

The third plot *Scale-Location plot* is used to test homoscedasticity assumption (constant variance across the sample)[4]. Again, the residuals behave erratically towards the extremes of the distribution. Two solutions to make our estimations consistent are *Huber-White standard errors* or *Boostrap standard errors* (not discussed in this report).

---

[4]The Barlett test is used to test statistically equal variance assumption.

*(a) Model 1*



*(b) Model 2*



*(c) Model 3*



*(d) Model 4*

**Figure C.1:** Diagnostic Plots for Linear Regression

19

# D. Discussing the Assumptions of ANOVA

## D.1. *Bartlett test*

Bartlett's test allows you to compare the variance of two or more samples to determine whether they are drawn from populations with equal variance. It is suitable for normally distributed data. The test has the null hypothesis that the variances are equal and the alterntive hypothesis that they are not equal.

**Data: Life expectancy by Continent**

**Bartlett's K-squared = 19.743, df = 5, p-value = 0.001397**

From the output we can see that the p-value of 0.001397 is less than the significance level of 0.05. This means we can reject the null hypothesis.



**Figure D.1: This plot shows the residuals (errors) on the y-axis and the fitted values (predicted values) on the x-axis. We can see that the points local in two group under 65 and over 72.5 and the residuals are not distributed equally above and blew zero. The red line is flat, then the relationship between the independent and dependent variable is linear.**

**Figure D.2: The Normal QQ plot indicates that the independent variable is approximately normal, since the observations approximately line-up with the Theoretically-derived normal value.**

**Figure D.3: The Scale-Location plot shows if residuals are spread equally along the ranges of predictors. This is how we can check the assumption of equal variance.**

## E. R code

```
1   ########################
2   ####### MA 317 #########
3   ##### Group Project ####
4   ## Clenaning the Data ##
5   ########################
6
7   #loading libraries
8   library(ggplot2)
9   library(pastecs)
10  library(psych)
11
12  #set wd
13  wd=setwd("D:/Documentos/Essex/Modelling Experimental Data/group project")
14
15  #reading the data set
16  data=read.csv("LifeExpectancy1.csv", header=T)
17
18  #asigning different names
19  names(data)=c("country", "code", "life_exp", "electricity", "adj_income", "primary_school",
    ↪  "exp_primary_ed", "ppp", "mortality_water", "literacy_rate", "pop_growth", "pop_total",
    ↪  "primary_completion", "secondary_education_duration", "sec_ed_teachers", "health_exp",
    ↪  "health_exp_pc", "unemployment_total", "unemployment_youth", "rural_pop", "adolscent_fert",
    ↪  "GDP_pc", "mobile_coverage", "internet_users")
20
21  #Creating a list of variables to be transformed to log form
22  lista=list("adj_income","primary_school","mortality_water","pop_total",
    ↪  "secondary_education_duration", "GDP_pc", "mobile_coverage", "sec_ed_teachers", "health_exp_pc")
23
24  #loop to iterate throught the list
25  for (l in lista){
26    print(l)
27    data[l]=log(data[l])
28  }
29
30  #save the data
31  write.csv(data, "LifeExpectancy1_log.csv", row.names = F)
32
33  ########################
34  ####### MA 317 #########
35  ##### Group Project ####
36  ########Task 1#########
37  #Descriptive Satatistics#
38  ########################
39  #model of life expectancy with all variables
40  model1 <- lm(life_exp ~ electricity, data=countries1)
41  model2 <- lm(life_exp ~ adj_income, data=countries1)
42  model3 <- lm(life_exp ~ primary_school, data=countries1)
43  model4 <- lm(life_exp ~ exp_primary_ed, data=countries1)
44  model5 <- lm(life_exp ~ ppp, data=countries1)
45  model6 <- lm(life_exp ~ mortality_water, data=countries1)
```

```
46   model7 <- lm(life_exp ~ literacy_rate, data=countries1)
47   model8 <- lm(life_exp ~ pop_growth, data=countries1)
48   model9 <- lm(life_exp ~ pop_total, data=countries1)
49   model10 <- lm(life_exp ~ primary_completion, data=countries1)
50   model11<- lm(life_exp ~ secondary_education_duration, data=countries1)
51   model12 <- lm(life_exp ~ sec_ed_teachers, data=countries1)
52   model13 <- lm(life_exp ~ health_exp, data=countries1)
53   model14 <- lm(life_exp ~ health_exp_pc, data=countries1)
54   model15 <- lm(life_exp ~ unemployment_total, data=countries1)
55   model16 <- lm(life_exp ~ unemployment_youth, data=countries1)
56   model17 <- lm(life_exp ~ rural_pop, data=countries1)
57   model18 <- lm(life_exp ~ adolscent_fert, data=countries1)
58   model19 <- lm(life_exp ~ GDP_pc, data=countries1)
59   model20 <- lm(life_exp ~ mobile_coverage, data=countries1)
60   model21 <- lm(life_exp ~ internet_users, data=countries1)
61   model19
62   #make text of r square for scatter plots
63   e1 <- paste("R^2 = ", round(summary(model1)$r.squared, 3))
64   e2 <- paste("R^2 = ", round(summary(model2)$r.squared, 3))
65   e3 <- paste("R^2 = ", round(summary(model3)$r.squared, 3))
66   e4 <- paste("R^2 = ", round(summary(model4)$r.squared, 3))
67   e5 <- paste("R^2 = ", round(summary(model5)$r.squared, 3))
68   e6 <- paste("R^2 = ", round(summary(model6)$r.squared, 3))
69   e7 <- paste("R^2 = ", round(summary(model7)$r.squared, 3))
70   e8 <- paste("R^2 = ", round(summary(model8)$r.squared, 3))
71   e9 <- paste("R^2 = ", round(summary(model9)$r.squared, 3))
72   e10 <- paste("R^2 = ", round(summary(model10)$r.squared, 3))
73   e11 <- paste("R^2 = ", round(summary(model11)$r.squared, 3))
74   e12 <- paste("R^2 = ", round(summary(model12)$r.squared, 3))
75   e13 <- paste("R^2 = ", round(summary(model13)$r.squared, 3))
76   e14 <- paste("R^2 = ", round(summary(model14)$r.squared, 3))
77   e15 <- paste("R^2 = ", round(summary(model15)$r.squared, 3))
78   e16 <- paste("R^2 = ", round(summary(model16)$r.squared, 3))
79   e17 <- paste("R^2 = ", round(summary(model17)$r.squared, 3))
80   e18 <- paste("R^2 = ", round(summary(model18)$r.squared, 3))
81   e19 <- paste("R^2 = ", round(summary(model19)$r.squared, 3))
82   e20 <- paste("R^2 = ", round(summary(model20)$r.squared, 3))
83   e21 <- paste("R^2 = ", round(summary(model21)$r.squared, 3))
84
85   #scatter plot with linear regression
86   gp1 = ggplot(countries1)+geom_point(aes(electricity, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(electricity, life_exp), method = "lm") +
     ↪  xlab("Electricity") + ylab("Life expectancy")+annotate(geom="text", x = 50, y= 80,
     ↪  label=e1,color="black", size= 3)
87   gp2 = ggplot(countries1)+geom_point(aes(adj_income, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(adj_income, life_exp), method = "lm") + xlab("Adjusted
     ↪  Income*") + ylab("Life expectancy")+annotate(geom="text", x=21, y=83, label=e2,color="black",
     ↪  size= 3)
```

```
88  gp3 = ggplot(countries1)+geom_point(aes(primary_school, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(primary_school, life_exp), method = "lm") +
    ↪  xlab("Children out of school*") + ylab("Life expectancy")+annotate(geom="text", x=1, y=90,
    ↪  label=e3,color="black", size= 3)
89  gp4 = ggplot(countries1)+geom_point(aes(exp_primary_ed, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(exp_primary_ed, life_exp), method = "lm") +
    ↪  xlab("Primary education expenditure")+ ylab("Life expectancy") +annotate(geom="text", x=50,
    ↪  y=80, label=e4,color="black", size= 3)
90  gp5 = ggplot(countries1)+geom_point(aes(ppp, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(ppp, life_exp), method = "lm") + xlab("PPP")+
    ↪  ylab("Life expectancy")+annotate(geom="text", x=2.5e+08, y=80, label=e4,color="black", size= 3)
91  gp6 = ggplot(countries1)+geom_point(aes(mortality_water, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(mortality_water, life_exp), method = "lm") +
    ↪  xlab("Unsafe water mortality*") + ylab("Life expectancy")+annotate(geom="text", x=3, y=80,
    ↪  label=e6,color="black", size= 3)
92  gp7 = ggplot(countries1)+geom_point(aes(literacy_rate, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(literacy_rate, life_exp), method = "lm") + xlab("Adult
    ↪  literacy rate") + ylab("Life expectancy")+annotate(geom="text", x=60, y=80,
    ↪  label=e7,color="black", size= 3)
93  gp8 = ggplot(countries1)+geom_point(aes(pop_growth, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(pop_growth, life_exp), method = "lm") +
    ↪  xlab("Population growth") + ylab("Life expectancy")+annotate(geom="text", x=3, y=87,
    ↪  label=e8,color="black", size= 3)
94  gp9 = ggplot(countries1)+geom_point(aes(pop_total, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(pop_total, life_exp), method = "lm") + xlab("Population
    ↪  total*") + ylab("Life expectancy") + xlim(11,21)+annotate(geom="text", x=20, y=82,
    ↪  label=e9,color="black", size= 3)
95  gp10 = ggplot(countries1)+geom_point(aes(primary_completion, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(primary_completion, life_exp), method = "lm") +
    ↪  xlab("Primary completion") + ylab("Life expectancy")+ xlim(40,120)+annotate(geom="text", x=60,
    ↪  y=79, label=e10,color="black", size= 3)
96  gp11 = ggplot(countries1)+geom_point(aes(secondary_education_duration, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(secondary_education_duration, life_exp), method = "lm")
    ↪  + xlab("Secondary education*") + ylab("Life expectancy")+annotate(geom="text", x=2.1, y=55,
    ↪  label=e11,color="black", size= 3)
97  gp12 = ggplot(countries1)+geom_point(aes(sec_ed_teachers, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(sec_ed_teachers, life_exp), method = "lm") +
    ↪  xlab("Secondary education teachers*") + ylab("Life expectancy")+annotate(geom="text", x=14,
    ↪  y=55, label=e12,color="black", size= 3) + xlim(6,15)
98  gp13 = ggplot(countries1)+geom_point(aes(health_exp, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(health_exp, life_exp), method = "lm") + xlab("Health
    ↪  expenditure") + ylab("Life expectancy")+ xlim(2,18)+annotate(geom="text", x=5, y=86,
    ↪  label=e13,color="black", size= 3)
99  gp14 = ggplot(countries1)+geom_point(aes(health_exp_pc, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(health_exp_pc, life_exp)) + xlab("Health expenditure
    ↪  per capita*") + ylab("Life expectancy")+annotate(geom="text", x=2500, y=95,
    ↪  label=e14,color="black", size= 3)
100 gp15 = ggplot(countries1)+geom_point(aes(unemployment_total, life_exp),
    ↪  colour="red")+theme_bw()+geom_smooth(aes(unemployment_total, life_exp), method = "lm") +
    ↪  xlab("Unemployment") + ylab("Life expectancy")+annotate(geom="text", x=12, y=62,
    ↪  label=e15,color="black", size= 3)
```

```r
101  gp16 = ggplot(countries1)+geom_point(aes(unemployment_youth, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(unemployment_youth, life_exp), method = "lm") +
     ↪  xlab("Youth Unemployment") + ylab("Life expectancy")+annotate(geom="text", x=25, y=60,
     ↪  label=e16,color="black", size= 3)
102  gp17 = ggplot(countries1)+geom_point(aes(rural_pop, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(rural_pop, life_exp), method = "lm") + xlab("Rural
     ↪  population") + ylab("Life expectancy")+annotate(geom="text", x=12, y=60,
     ↪  label=e17,color="black", size= 3)
103  gp18 = ggplot(countries1)+geom_point(aes(adolscent_fert, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(adolscent_fert, life_exp), method = "lm") +
     ↪  xlab("Adolescent fertility rate") + ylab("Life expectancy")+annotate(geom="text", x=120, y=80,
     ↪  label=e18,color="black", size= 3)
104  gp19 = ggplot(countries1)+geom_point(aes(GDP_pc, life_exp), colour="red", size =
     ↪  3)+theme_bw()+theme(axis.title=element_text(size=15))+geom_smooth(aes(GDP_pc, life_exp), method
     ↪  = "lm") + xlab("GDP per capita*") + ylab("Life expectancy")+annotate(geom="text", x=8, y=83,
     ↪  label=e19,color="black", size= 7)
105  gp20 = ggplot(countries1)+geom_point(aes(mobile_coverage, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(mobile_coverage, life_exp), method = "lm") +
     ↪  xlab("Mobile subscriptions*") + ylab("Life expectancy")+annotate(geom="text", x=3.2, y=82,
     ↪  label=e20,color="black", size= 3)
106  gp21 = ggplot(countries1)+geom_point(aes(internet_users, life_exp),
     ↪  colour="red")+theme_bw()+geom_smooth(aes(internet_users, life_exp), method = "lm") +
     ↪  xlab("Internet users") + ylab("Life expectancy")+annotate(geom="text", x=25, y=82,
     ↪  label=e21,color="black", size= 3)
107  #appendix
108  plg1 =  plot_grid(gp1, gp2,gp3,gp4,gp5,gp6,gp7,gp8,gp9, ncol=3)
109  #show output 6 scatterplot
110  plg2 =  plot_grid(gp10, gp11,gp12,gp13,gp14,gp15,gp16,gp17,gp18, ncol=3)
111  #show output 6 scatterplot
112  plg3 =  plot_grid(gp20, gp21, ncol=3, nrow=3)
113  #show output 2 scatterplot
114
115  options(scipen=100)
116  options(digits=3)
117  head(countries)
118  #Des. stat question 1
119  #Adjusted Income in billion
120  countries$adj_income = countries$adj_income/1000000000
121  #PPP in million
122  countries$ppp = countries$ppp/1000000
123  #population total in million
124  countries$pop_total = countries$pop_total/1000000
125  #secondary education teachers in thousand
126  countries$sec_ed_teachers = countries$sec_ed_teachers/1000
127  #exclude Country and Code column
128  table=stat.desc(countries[,-c(1,2)], basic=F)
129  #transpose the table
130  table=t(table[-5,])
131  #extract only SD column
132  SD = data.frame(table[,5])
133  #summary for every variable in transpose
```

```r
134  stattable = t(summary(countries$life_exp))
135  stattable1 = t(summary(countries$electricity))
136  stattable2 = t(summary(countries$adj_income))
137  stattable3 = t(summary(countries$primary_school))
138  stattable4 = t(summary(countries$exp_primary_ed))
139  stattable5 = t(summary(countries$ppp))
140  stattable6 = t(summary(countries$mortality_water))
141  stattable7 = t(summary(countries$literacy_rate))
142  stattable8 = t(summary(countries$pop_growth))
143  stattable9 = t(summary(countries$pop_total))
144  stattable10 = t(summary(countries$primary_completion))
145  stattable11 = t(summary(countries$secondary_education_duration))
146  stattable12 = t(summary(countries$sec_ed_teachers))
147  stattable13 = t(summary(countries$health_exp))
148  stattable14 = t(summary(countries$health_exp_pc))
149  stattable15 = t(summary(countries$unemployment_total))
150  stattable16 = t(summary(countries$unemployment_youth))
151  stattable17 = t(summary(countries$rural_pop))
152  stattable18 = t(summary(countries$adolscent_fert))
153  stattable19 = t(summary(countries$GDP_pc))
154  stattable20 = t(summary(countries$mobile_coverage))
155  stattable21 = t(summary(countries$internet_users))
156  #row combine for all summary
157  rb = rbind(stattable, stattable1,stattable2,stattable3, stattable4,stattable5,stattable6,
158  stattable7,stattable8, stattable9,stattable10, stattable11,stattable12,stattable13,stattable14,
159  stattable15,stattable16, stattable17, stattable18, stattable19,stattable20,stattable21)
160  #make it to be data frame
161  rbdf = data.frame(rb)
162  #column combine with Standard Deviation
163  Dstattable = cbind(SD,rbdf)
164  #delete column 1st.qu and 3rd.qu
165  DFdes = Dstattable[,-c(3,6)]
166  #rename the columns
167  colnames(DFdes) = c("Standard Deviation","Min.","Median","Mean","Max.","NA.s")
168  #reorder the columns
169  DFdes1 = DFdes[,c(4,3,1,2,5,6)]
170  #des.stat of all variables
171  DFdes1
172  #make the table to the latex format
173  kable(DFdes1, format="latex", digits=2, booktabs=TRUE)
174  #variables which have most significant r square
175  countries2 = select(countries, c(life_exp,electricity,literacy_rate, primary_completion,health_exp,
     ↪  health_exp_pc,GDP_pc,
176  mobile_coverage,internet_users, exp_primary_ed,mortality_water, rural_pop,pop_growth,adj_income,
177  adolscent_fert))
178  #geom boxplot
179  attach(countries2)
180  bb1 = ggplot(countries2, aes(x= "",y=life_exp))+geom_boxplot(fill = "#FF6347")+ylab("Life
     ↪  expectancy")+xlab("")+
181  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
```

```
182  bb2 = ggplot(countries2, aes(x= "",y=electricity))+geom_boxplot(fill =
     ↪  "#FF6347")+ylab("Electricity")+xlab("")+
183  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
184  bb3 = ggplot(countries2, aes(x= "",y=literacy_rate))+geom_boxplot(fill = "#FF6347")+ylab("Adult
     ↪  literacy rate")+xlab("")+
185  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
186  bb4 = ggplot(countries2, aes(x= "",y=primary_completion))+geom_boxplot(fill =
     ↪  "#FF6347")+ylab("Primary completion")+xlab("")+
187  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
188  bb5 = ggplot(countries2, aes(x= "",y=health_exp))+geom_boxplot(fill = "#FF6347")+ylab("Health
     ↪  expenditure")+xlab("")+
189  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
190  bb6 = ggplot(countries2, aes(x= "",y=health_exp_pc))+geom_boxplot(fill = "#FF6347")+ylab("Health
     ↪  expenditure per capita")+xlab("")+
191  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
192  bb7 = ggplot(countries2, aes(x= "",y=GDP_pc))+geom_boxplot(fill = "#FF6347" , width = .3)+ylab("GDP
     ↪  per capita")+xlab("")+
193  stat_boxplot(geom ='errorbar', width = .3) + coord_flip() + theme_bw()
194  bb8 = ggplot(countries2, aes(x= "",y=mobile_coverage))+geom_boxplot(fill = "#FF6347")+ylab("Mobile
     ↪  subscriptions")+xlab("")+
195  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
196  bb9 = ggplot(countries2, aes(x= "",y=internet_users))+geom_boxplot(fill = "#FF6347")+ylab("Internet
     ↪  users")+xlab("")+
197  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
198  bb10 = ggplot(countries2, aes(x= "",y=exp_primary_ed))+geom_boxplot(fill = "#FF6347")+ylab("Primary
     ↪  education expenditure")+xlab("")+
199  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
200  bb11 = ggplot(countries2, aes(x= "",y=mortality_water))+geom_boxplot(fill = "#FF6347")+ylab("Unsafe
     ↪  water mortality")+xlab("")+
201  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
202  bb12 = ggplot(countries2, aes(x= "",y=rural_pop))+geom_boxplot(fill = "#FF6347")+ylab("Rural
     ↪  population")+xlab("")+
203  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
204  bb13 = ggplot(countries2, aes(x= "",y=pop_growth))+geom_boxplot(fill = "#FF6347")+ylab("Population
     ↪  Growth")+xlab("")+
205  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
206  bb14 = ggplot(countries2, aes(x= "",y=adolscent_fert))+geom_boxplot(fill =
     ↪  "#FF6347")+ylab("Adolscent fertility rate")+xlab("")+
207  stat_boxplot(geom ='errorbar') + coord_flip() + theme_bw()
208  #variables
209  bbmain = plot_grid(bb1,bb2, bb11,bb13,bb4,bb6,bb12,bb14,bb9,bb8, ncol=2)
210
211  #######################
212  ####### MA 317 #########
213  ##### Group Project ####
214  ####### Task 2 #########
215  ##### Imputation#####
216  #######################
217  library(ggpubr)
218  library(ggplot2)
219  library(grid)
```

26

```
220  library(gridExtra)
221  library(cowplot)
222  setwd('/Users/preciousakinyele/Downloads')
223  data=read.csv("LifeExpectancy1_log.csv", header=T)
224  head(data)
225  countries=data[1:217,]
226
227  #USING MEAN IMPUTATION
228  #Mean Imputation for Log(Adjusted Income)
229  countries$adj_income[is.na(countries$adj_income)]<- mean(countries$adj_income, na.rm=TRUE)
230
231  #Mean Imputation for Access to Electricity
232  countries$electricity[is.na(countries$electricity)]<- mean(countries$electricity, na.rm = TRUE)
233
234  #Mean Imputation for Population growth
235  countries$pop_growth[is.na(countries$pop_growth)]<- mean(countries$pop_growth, na.rm = TRUE)
236
237  #Mean Imputation for Log(Population total)
238  countries$pop_total[is.na(countries$pop_total)]<- mean(countries$pop_total, na.rm = TRUE)
239
240  #Mean Imputation for Primary Completion
241  countries$primary_completion[is.na(countries$primary_completion)]<-mean(countries$primary_completion,
     ↪  na.rm = TRUE)
242
243  #Mean Imputation for Log(Secondary Education Duration)
244  countries$secondary_education_duration[is.na(countries$secondary_education_duration)]<-
     ↪  mean(countries$secondary_education_duration, na.rm = TRUE)
245
246  #Mean Imputation for Health Expenditure
247  countries$health_exp[is.na(countries$health_exp)]<-mean(countries$health_exp, na.rm = TRUE)
248
249  #Mean Imputation for Rural Population
250  countries$rural_pop[is.na(countries$rural_pop)]<- mean(countries$rural_pop, na.rm = TRUE)
251
252  #Mean Imputation for Log(Mobile Cellular Subscription)
253  countries$mobile_coverage[is.na(countries$mobile_coverage)]<- mean(countries$mobile_coverage, na.rm
     ↪  = TRUE)
254
255  #Mean Imputation for Internet Users
256  countries$internet_users[is.na(countries$internet_users)]<- mean(countries$internet_users, na.rm =
     ↪  TRUE)
257
258  #USING MEDIAN IMPUTATION DUE TO A LARGE DIFFERENCE IN MEAN AND MEDIAN. RESULTING FROM OUTLIERS
259  #Median imputation for Children out of Primary school
260  sum.primary=summary((countries$primary_school))[3]
261  sum.primary
262  countries$primary_school[is.na(countries$primary_school)]<- sum.primary
263
264  #Median imputation for Log(GDP per capita)
265  sum.gdp=summary((countries$GDP_pc))[3]
266  sum.gdp
```

```r
267  countries$GDP_pc[is.na(countries$GDP_pc)]<- sum.gdp
268
269  #Median imputation for Log(Health expenditure per capita)
270  sum.health_pc=summary((countries$health_exp_pc))[3]
271  sum.health_pc
272  countries$health_exp_pc[is.na(countries$health_exp_pc)]<- sum.health_pc
273
274  #Median imputation for Log(Mortality rate due to unsafe water)
275  sum.mortarlity_water=summary((countries$mortality_water))[3]
276  sum.mortarlity_water
277  countries$mortality_water[is.na(countries$mortality_water)]<- sum.mortarlity_water
278
279  #Median imputation for Adolescent Fertility rate
280  sum.adolscent_fert=summary((countries$adolscent_fert))[3]
281  sum.adolscent_fert
282  countries$adolscent_fert[is.na(countries$adolscent_fert)]<- sum.adolscent_fert
283
284  #LINEAR REGRESSION METHOD FOR DETERMINISTIC IMPUTATION
285  #Regression method for Primary education expenditure,
286  lm_exp<-
     ↪  (lm(countries$exp_primary_ed~countries$electricity+countries$adj_income+countries$primary_school+countries$mo
287  countries$health_exp_pc+countries$rural_pop+countries$adolscent_fert+countries$GDP_pc+countries$mobile_coverage+c
288  summary(lm_exp)
289  pred1 <- predict(lm_exp)
290  impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
291  countries$exp_primary_ed<-impute(countries$exp_primary_ed,pred1)
292
293  #Regression method for Adult literacy rate
294  #LITERACY RATE HAS A LOT OF MISSING VALUES, NEVERTHELESS, HERE IS THE LINEAR REGRESSION FOR IT
295  lm_litrate<-(lm(countries$literacy_rate~countries$electricity+countries$adj_income+countries$primary_school+count
296  pred2 <- predict(lm_litrate)
297  countries$literacy_rate <-impute(countries$literacy_rate,pred2)
298
299
300  #Regression method for Log(Secondary education teachers)
301  lm_seced<-(lm(countries$sec_ed_teachers~countries$electricity+countries$adj_income+countries$primary_school+count
302  pred3 <- predict(lm_seced)
303  countries$sec_ed_teachers<-impute(countries$sec_ed_teachers,pred3)
304
305
306  #Regression method for Total Unemployment
307  lm_total.unemp<-(lm(countries$unemployment_total~countries$electricity+countries$adj_income+countries$primary_sch
308  pred4 <- predict(lm_total.unemp)
309  countries$unemployment_total<-impute(countries$unemployment_total,pred4)
310
311  #Regression method Total Youth Unemployment
312  lm_total.youth<-(lm(countries$unemployment_youth~countries$electricity+countries$adj_income+countries$primary_sch
313  pred5 <- predict(lm_total.youth)
314  countries$unemployment_youth <- impute(countries$unemployment_youth,pred5)
315
316  #creating a csv file with imputations
```

```
317   write.csv(countries, "imputation.csv", row.names = F)

318

319   #PDF PLOTS
320   #Plotting the Probability Density Function for Log (Adjusted Income) estimator variable
321   countries$a = countries$adj_income
322   countries$adj_income[is.na(countries$adj_income)]<- mean(countries$adj_income, na.rm=TRUE)
323   gp1 = ggplot()+geom_line(aes(x=na.omit(countries$adj_income) , y= stat(density)), stat =
      ↪  'density')+theme_bw()+theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+ geom_line(aes(countries$a),stat = 'density', color='red')+xlab("Adjusted Income*")

324

325   #Plotting the Probability Density Function for Access to Electricity estimator variable
326   countries$b= countries$electricity
327   countries$electricity[is.na(countries$electricity)]<- mean(countries$electricity, na.rm = TRUE)
328   gp2 = ggplot()+geom_line(aes(x=na.omit(countries$electricity), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
329   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$b),stat = 'density', color='red')+xlab("Electricity")

330

331   #Plotting the Probability Density Function for Log Population growth estimator variable
332   countries$c = countries$pop_growth
333   countries$pop_growth[is.na(countries$pop_growth)]<- mean(countries$pop_growth, na.rm = TRUE)
334   gp3 = ggplot()+geom_line(aes(x=na.omit(countries$pop_growth), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
335   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$c),stat = 'density', color='red')+ xlab("Population growth")

336

337   #Plotting the Probability Density Function for Log (Total population) estimator variable
338   countries$d = countries$pop_total
339   countries$pop_total[is.na(countries$pop_total)]<- mean(countries$pop_total, na.rm = TRUE)
340   gp4 = ggplot()+geom_line(aes(x=na.omit(countries$pop_total), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
341   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$d),stat = 'density', color='red')+xlab("Population total*")

342

343   #Plotting the Probability Density Function for Primary Completion estimator variable
344   countries$e = countries$primary_completion
345   countries$primary_completion[is.na(countries$primary_completion)]<-mean(countries$primary_completion,
      ↪  na.rm = TRUE)
346   gp5 = ggplot()+geom_line(aes(x=na.omit(countries$primary_completion), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
347   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$e),stat = 'density', color='red')+xlab("Primary completion")

348

349   #Plotting the Probability Density Function for Log (Secondary Education) estimator variable
350   countries$f = countries$secondary_education_duration
351   countries$secondary_education_duration[is.na(countries$secondary_education_duration)]<-
      ↪  mean(countries$secondary_education_duration, na.rm = TRUE)
352   gp6 = ggplot()+geom_line(aes(x=na.omit(countries$secondary_education_duration), y= stat(density)),
      ↪  stat ='density')+theme_bw()+
353   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$f),stat = 'density', color='red')+xlab("Secondary education*")
```

```r
354
355   #Plotting the Probability Density Function for Health Expenditure estimator variable
356   countries$g = countries$health_exp
357   countries$health_exp[is.na(countries$health_exp)]<-mean(countries$health_exp, na.rm = TRUE)
358   gp7 = ggplot()+geom_line(aes(x=na.omit(countries$health_exp), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
359   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$g),stat = 'density', color='red')+xlab("Health expenditure")
360
361   #Plotting the Probability Density Function for Rural Population estimator variable
362   countries$h = countries$rural_pop
363   countries$rural_pop[is.na(countries$rural_pop)]<- mean(countries$rural_pop, na.rm = TRUE)
364   gp8 = ggplot()+geom_line(aes(x=na.omit(countries$rural_pop), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
365   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$h),stat = 'density', color='red')+xlab("Rural population")
366
367   #Plotting the Probability Density Function for Log (Mobile subscription) estimator variable denoted
      ↪  as Mobile Subscription*
368   countries$i = countries$mobile_coverage
369   countries$mobile_coverage[is.na(countries$mobile_coverage)]<- mean(countries$mobile_coverage, na.rm
      ↪  = TRUE)
370   gp9 = ggplot()+geom_line(aes(x=na.omit(countries$mobile_coverage), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
371   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$i),stat = 'density', color='red')+xlab("Mobile subscription*")
372
373   #Plotting the Probability Density Function for Internet Users estimator variable
374   countries$j = countries$internet_users
375   countries$internet_users[is.na(countries$internet_users)]<- mean(countries$internet_users, na.rm =
      ↪  TRUE)
376   gp10 = ggplot()+geom_line(aes(x=na.omit(countries$internet_users), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
377   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$j),stat = 'density', color='red')+xlab("Internet users")
378
379   #Plotting the Probability Density Function for Children out of primary school estimator variable
380   countries$k = countries$primary_school
381   sum.primary=summary((countries$primary_school))[3]
382   countries$primary_school[is.na(countries$primary_school)]<- sum.primary
383   gp11 = ggplot()+geom_line(aes(x=na.omit(countries$primary_school), y= stat(density)), stat =
      ↪  'density')+theme_bw()+
384   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪  9))+geom_line(aes(countries$k),stat = 'density', color='red')+xlab("Children out of school")
385
386   #Plotting the Probability Density Function for Log (GDP per capita) estimator variable denoted as GP
      ↪  per capita*
387   countries$l = countries$GDP_pc
388   sum.gdp=summary((countries$GDP_pc))[3]
389   countries$GDP_pc[is.na(countries$GDP_pc)]<- sum.gdp
```

```r
390  gp12 = ggplot()+geom_line(aes(x=na.omit(countries$GDP_pc), y= stat(density)), stat =
     ↪  'density')+theme_bw()+
391  theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
     ↪  9))+geom_line(aes(countries$l),stat = 'density', color='red')+xlab("GDP per capita*")
392
393  #Plotting the Probability Density Function for Log (Health expenditure per capita) estimator
     ↪  variable denoted as Health expenditure per capita*
394  countries$m = countries$health_exp_pc
395  sum.health_pc=summary((countries$health_exp_pc))[3]
396  countries$health_exp_pc[is.na(countries$health_exp_pc)]<- sum.health_pc
397  gp13 = ggplot()+geom_line(aes(x=na.omit(countries$health_exp_pc), y= stat(density)), stat =
     ↪  'density')+theme_bw()+
398  theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
     ↪  9))+geom_line(aes(countries$m),stat = 'density', color='red')+ xlab("Health expenditure per
     ↪  capita*")
399
400  #Plotting the Probability Density Function for Log (Mortality rate due to unsafe water) estimator
     ↪  variable denoted as Unsafe water mortality*
401  countries$n = countries$mortality_water
402  sum.mortarlity_water=summary((countries$mortality_water))[3]
403  countries$mortality_water[is.na(countries$mortality_water)]<- sum.mortarlity_water
404  gp14 = ggplot()+geom_line(aes(x=na.omit(countries$mortality_water), y= stat(density)), stat =
     ↪  'density')+theme_bw()+
405  theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
     ↪  9))+geom_line(aes(countries$n),stat = 'density', color='red')+xlab("Unsafe water mortality*")
406
407  #Plotting the Probability Density Function for Adolescent Fertility rate estimator variable
408  countries$o = countries$adolscent_fert
409  sum.adolscent_fert=summary((countries$adolscent_fert))[3]
410  countries$adolscent_fert[is.na(countries$adolscent_fert)]<- sum.adolscent_fert
411  gp15 = ggplot()+geom_line(aes(x=na.omit(countries$adolscent_fert), y= stat(density)), stat =
     ↪  'density')+theme_bw()+
412  theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
     ↪  9))+geom_line(aes(countries$o),stat = 'density', color='red')+xlab("Adolescent fertilty rate ")
413
414  #Plotting the Probability Density Function for Primary education expenditure estimator variable
415  countries$p = countries$exp_primary_ed
416  lm_exp<-
     ↪  (lm(countries$exp_primary_ed~countries$electricity+countries$adj_income+countries$primary_school+countries$mo
417  pred1 <- predict(lm_exp)
418  impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
419  countries$exp_primary_ed <-+ impute(countries$exp_primary_ed,pred1)
420  gp16 = ggplot()+geom_line(aes(x=na.omit(countries$exp_primary_ed), y= stat(density)), stat =
     ↪  'density')+theme_bw()+
421  theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
     ↪  9))+geom_line(aes(countries$p),stat = 'density', color='red')+xlab("Primary education
     ↪  expenditure")
422
423  #Plotting the Probability Density Function for Adult literacy rate estimator variable
424  countries$q = countries$literacy_rate
425  lm_litrate<-(lm(countries$literacy_rate~countries$electricity+countries$adj_income+countries$primary_school+count
```

```r
426   pred2 <- predict(lm_litrate)
427   impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
428   countries$literacy_rate<-+impute(countries$literacy_rate,pred2)
429
430   gp17 = ggplot()+geom_line(aes(x=na.omit(countries$literacy_rate), y= stat(density)), stat =
      ↪   'density')+theme_bw()+
431   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪   9))+geom_line(aes(countries$q),stat = 'density', color='red')+xlab("Adult literacy rate")
432
433   #Plotting the Probability Density Function for Log (secondary education teachers) estimator variable
      ↪   denoted as Secondary education teachers*
434   countries$r = countries$sec_ed_teachers
435   lm_seced<-(lm(countries$sec_ed_teachers~countries$electricity+countries$adj_income+countries$primary_school+count
436   pred3 <- predict(lm_seced)
437   impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
438   countries$sec_ed_teachers<-+impute(countries$sec_ed_teachers,pred3)
439   gp18 = ggplot()+geom_line(aes(x=na.omit(countries$sec_ed_teachers), y= stat(density)), stat =
      ↪   'density')+theme_bw()+
440   theme(axis.title.y = element_blank(),axis.title.x = element_text(size =
      ↪   9))+geom_line(aes(countries$r),stat = 'density', color='red')+xlab("Secondary education
      ↪   teachers*")
441
442   #Plotting the Probability Density Function for Total Unemployment estimator variable
443   countries$s = countries$unemployment_total
444   lm_total.unemp<-(lm(countries$unemployment_total~countries$electricity+countries$adj_income+countries$primary_sch
445   pred4 <- predict(lm_total.unemp)
446   impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
447   countries$unemployment_total<-+impute(countries$unemployment_total,pred4)
448   gp19 = ggplot()+geom_line(aes(x=na.omit(countries$unemployment_total), y= stat(density)), stat =
      ↪   'density')+theme_bw()+
449   theme(axis.title.y = element_blank(),axis.title.x = element_text(size = 9))+
      ↪   geom_line(aes(countries$s),stat = 'density', color='red')+xlab("Unemployment")
450
451   #Plotting the Probability Density Function for Total Youth Unemployment estimator variable
452   countries$t = countries$unemployment_youth
453   lm_youth.total<-(lm(countries$unemployment_youth~countries$electricity+countries$adj_income+countries$primary_sch
454   pred5 <- predict(lm_youth.total)
455   impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
456   countries$unemployment_youth <-+ impute(countries$unemployment_youth,pred5)
457   gp20 = ggplot()+geom_line(aes(x=na.omit(countries$unemployment_youth), y= stat(density)), stat =
      ↪   'density')+theme_bw()+
458   theme(axis.title.y = element_blank(),axis.title.x = element_text(size = 9)) +
      ↪   geom_line(aes(countries$t),stat = 'density', color='red')+xlab("Youth unemployment")
459   gp20
460
461   #A PDF plot of all the estimator variables from gp1 to gp9 agaisnt density in grids
462   rplots1= plot_grid(gp1, gp2, gp3, gp4, gp5, gp6, gp7, gp8, gp9, ncol = 3)
463   rplots1
464
465   #A PDF plot of all the estimator variables from gp10 to gp18 agaisnt density in grids
466   rplots2= plot_grid(gp10, gp11, gp12, gp13, gp14, gp15, gp16, gp17, gp18, ncol=3)
```

```
467   rplots2

468

469   #A PDF plot of all the estimator variables from gp19 to gp20 agaisnt density in grids
470   rplots3= plot_grid(gp19, gp20, ncol=3, nrow =3)
471   rplots3

472

473   #MISSING DATA IMPUTATION METHODS JUSTIFICATION

474

475   #reload the csv file for justification
476   data=read.csv("LifeExpectancy1_log.csv", header=T)

477

478   #To check if there are any countries with complete estimator variables.
479   countries %>% filter(complete.cases(.))

480

481   #to check the countries with life expectancies before carrying out the complete case
482   countries$life_exp #211

483

484   #Since there are no countries with complete estimator variables, therefore, conditions for complete
485   #cases were relaxed. We then assumed a complete case by
486   #removing ppp, literacy rate and Primary education expenditure
487   b2= countries %>% filter(c(!is.na(countries$electricity)& !is.na(countries$adj_income)&
      ↪   !is.na(countries$primary_school)& !is.na(countries$pop_growth)& !is.na(countries$pop_total)&
488   !is.na(countries$primary_completion)&!is.na(countries$secondary_education_duration)&!is.na(countries$sec_ed_teach
      ↪   !is.na(countries$adolscent_fert)&!is.na(countries$unemployment_youth)&!is.na(countries$mortality_water)&!is.n
      ↪   !is.na(countries$mobile_coverage)&!is.na(countries$internet_users)))

489

490   #to check number of countries that have the assumed complete case.
491   count(b2)

492

493   #to check number of countries with life expectancy under the assumed complete case.
494   b2$life_exp #48

495

496   #to plot the graph of life expectancy under the assumed complete case, as well as life expectancy
497   #having all the predictor variables.
498   ggplot()+geom_line(data=b2,aes(x= b2$life_exp), stat = 'density')+theme_bw()+
499   geom_line(data=countries,aes(x= countries$life_exp), stat = 'density', color='red')+xlab('Life
      ↪   Expectancy')
500   aspect_ratio<-2
501   height <-7
502   ggsave( height = 7, width = 7*aspect_ratio, "completecase.pdf")

503

504

505   #######################
506   ####### MA 317 #########
507   ##### Group Project ####
508   ####### Task 4 #########
509   ##### Linear Model#####
510   #######################

511

512   #Cleaning workingspace
513   rm(list=ls())
```

```
514
515    #loading libraries
516    library(data.table)
517    library(ggplot2)
518    library(ggpubr)
519    library(stargazer)
520    library(quantreg)
521    library(kableExtra)
522

523

524    #Setting wd
525    wd=setwd("D:/Documentos/Essex/Modelling Experimental Data/Group Project")
526    countries=read.csv("imputation.csv", header=T) #Read file
527    head(countries)    #check if it's the correct data
528    attach(countries) #attach variables to ws
529

530    #4. Suggest a model which explains life expectancies in 2016. Justify you answer. Can this model be
     ↪   used
531    #to predict life expectancies of other countries which have not provided in 2016 data on life
     ↪   expectancy?
532    #adding a variable GINI (Income inequality)
533    gini=fread("WDIData.csv", header=T)   #read data
534    gini=gini[gini$`Indicator Name`=="GINI index (World Bank estimate)",c(1,2,61)]  #keeping just year
     ↪   of interest
535

536    countries=merge(countries, gini[,c(1,3)], by.x="country", by.y="Country Name")   #merging with
     ↪   country data
537    names(countries)[names(countries)=="2016"]="gini"
538    attach(countries) #attach variables to ws
539    rm(gini)
540

541    #we analyse the gini data
542    summary(gini)
543

544    #perform regression imputation
545    model_gini=summary(lm(gini~countries$electricity+countries$adj_income+
     ↪   countries$primary_school+countries$mortality_water+countries$pop_growth+
     ↪   countries$pop_total+ countries$primary_completion+countries$secondary_education_duration+
     ↪   countries$health_exp+countries$health_exp_pc+ countries$rural_pop+countries$adolscent_fert+
     ↪   countries$GDP_pc+countries$mobile_coverage+countries$internet_users))
546

547    X=matrix(1,nrow(countries))
548

549    X=cbind(X, countries$electricity,countries$adj_income,
     ↪   countries$primary_school,countries$mortality_water,countries$pop_growth, countries$pop_total,
     ↪   countries$primary_completion, countries$secondary_education_duration, countries$health_exp,
     ↪   countries$health_exp_pc, countries$rural_pop,countries$adolscent_fert, countries$GDP_pc,
     ↪   countries$mobile_coverage, countries$internet_users)
550

551    countries$gini_i=X%*%as.matrix(model_gini$coefficients[,1])
552
```

```
553   #impute <- function (a, a.impute){ ifelse (is.na(a), a.impute, a)}
554   #countries$gini_i <- impute(countries$gini,model_gini)
555   attach(countries) #attach variables to ws
556
557   #summarizing the new imputed var
558   summary(countries$gini_i)
559   #comparing density function before &after imputation
560   countries$gini_i=ifelse(countries$gini_i<25,mean(countries$gini_i),countries$gini_i) #min value
      ↪  before imp
561
562   ggplot()+geom_line(data=countries, aes(gini), stat="density")+theme_bw()+
      ↪  geom_line(data=as.data.frame(gini_i), aes(gini_i), stat="density", color="red")
563
564   #plotting correlation with life expectancy
565   gp1 = ggplot(countries)+geom_point(aes(log(gini_i), life_exp), colour="red")+
      ↪  theme_bw()+geom_smooth(aes(log(gini_i), life_exp), method = "lm") + xlab("log(Gini)") +
      ↪  ylab("Life expectancy")+stat_cor(aes(log(gini_i), life_exp,label = ..r.label..),method =
      ↪  "pearson", label.x = 3.97, label.y = 80)
566
567   gp1
568   ggsave("gini.pdf")
569   countries$gini=NULL
570
571   #Defining the saturated model with only those variables with high correlation (meausred
572   #by the correlation coefficient) to Life Expectancy
573       #Y: Life Expectancy
574       #matrix X: Electricity, water, mortality,  population growth, adjusted income
575       #primary completion rate, health expenditure pc
576       #and squared (see scatterplot), rural population, adolescent fertilty rate, internet users,
577       #mobile coverage
578
579   sub=countries[,-c(1,2,8,26)]
580   corr <- round(cor(na.omit(sub)), 1)
581   corr=melt(corr)
582
583   n=rev(list("Life Expectancy","Electricity","Adjusted Income", "Primary School", "Expenditure
      ↪  Primary","Water Mortality","Literacy Rate", "Population Growth", "Total Population", "Primary
      ↪  Completion", "Secondary Education", "Secondary Teachers", "Health Exp.","Health Exp. pc",
      ↪  "Unemployment", "Youth Unemploymnent", "Rural Opulation", "Adolescent Fertility", "GDP pc",
      ↪  "Mobile Coverage", "Internet Users", "Gini"))
584
585   ggplot(data = corr, aes(x=Var1,ordered(Var2, levels =    rev(sort(unique(Var2)))), fill=value)) +
586   geom_tile()+ theme(axis.title.x = element_blank(),axis.text.y = element_text(vjust=0,
      ↪  axis.title.y = element_blank(),axis.text.x=element_blank(), legend.title = element_blank()) +
      ↪  scale_fill_gradient2(low="red", mid="yellow", high="red",
      ↪  midpoint=0,labels=c("-1.0", "-0.5","0","0.5","1.0"),breaks=c(-1,-.5,0,.5,1),limits=c(-1,1))+
587   scale_y_discrete(labels=n) ggsave("matrix_corr.pdf")
588
589   #From section 1 we know there is a non-linear correlation between life expectancy and
590   # health exp pc in levels. We add an extra squared term to capture this relation.
591
```

```
592  countries$health_pc_levels=exp(health_exp_pc)/1000
593  attach(countries)
594
595  #We perform stepwise variable selection with both AIC and BIC
596
597  best1=step(lm(life_exp~electricity +mortality_water+ pop_growth+primary_completion+health_pc_levels
     ↪ +
598  I(health_pc_levels^2)+ rural_pop+adolscent_fert +internet_users+mobile_coverage+adj_income),
     ↪ direction = "both")
599
600  aic1=best1$anova$AIC[length(best1$anova$AIC)]
601  extractAIC(best1, scale=0, k=log(nrow(countries)))
602
603  best1_1=step(lm(life_exp~electricity+ mortality_water+
     ↪ pop_growth+primary_completion+health_pc_levels +
604  I(health_pc_levels^2)+ rural_pop +adolscent_fert+internet_users+mobile_coverage+adj_income),
     ↪ direction = "both", k = log(nrow(countries)))
605
606  bic1=best1_1$anova$AIC[length(best1_1$anova$AIC)]
607
608  best2=step(lm(life_exp~electricity+ pop_growth+primary_completion+health_pc_levels+
     ↪ I(health_pc_levels^2)+rural_pop+adolscent_fert+ internet_users+mobile_coverage+ adj_income),
     ↪ direction = "both")
609
610  aic2=best2$anova$AIC[length(best2$anova$AIC)]
611  extractAIC(best2, scale=0, k=log(nrow(countries)))
612
613  best2_2=step(lm(life_exp~electricity+ pop_growth+primary_completion+health_pc_levels+
614  I(health_pc_levels^2)+ rural_pop+adolscent_fert+ internet_users+mobile_coverage+adj_income),
     ↪ direction = "both", k = log(nrow(countries)))
615
616  bic2=best2_2$anova$AIC[length(best2_2$anova$AIC)]
617  extractAIC(best2_2, scale=0)
618
619  ########################
620  #Summarize best models
621  #####################
622
623  diagPlot<-function(model){
624    p1<-ggplot(model, aes(.fitted, .resid))+geom_point()
625    p1<-p1+stat_smooth(method="loess")+geom_hline(yintercept=0, col="red", linetype="dashed")
626    p1<-p1+xlab("Fitted values")+ylab("Residuals")
627    p1<-p1+ggtitle("Residual vs Fitted Plot")+theme_bw()+theme(title = element_text(size=9))
628
629    p2<-ggplot(model, aes(qqnorm(.stdresid)[[1]], .stdresid))+geom_point(na.rm = TRUE)
630    p2<-p2+geom_abline()+xlab("Theoretical Quantiles")+ylab("Standardized Residuals")
631    p2<-p2+ggtitle("Normal Q-Q")+theme_bw()+theme(title = element_text(size=9))
632
633    p3<-ggplot(model, aes(.fitted, sqrt(abs(.stdresid))))+geom_point(na.rm=TRUE)
634    p3<-p3+stat_smooth(method="loess", na.rm = TRUE)+xlab("Fitted Value")
635    p3<-p3+ylab(expression(sqrt("|Standardized residuals|")))
```

```r
636    p3<-p3+ggtitle("Scale-Location")+theme_bw()+theme(title = element_text(size=9))
637    p3
638    return(ggarrange(p1, p2, p3, ncol=3, nrow=2))
639  }
640
641  #############
642  #####Model 1
643  #############
644  summary(best1)
645
646  #diagnostic plots
647  diagPlot(best1)
648  ggsave("best1.pdf")
649
650  #############
651  #####Model 2
652  #############
653  summary(best1_1)
654
655  #diagnostic plots
656  diagPlot(best1_1)
657  ggsave("best1_1.pdf")
658
659  #############
660  #####Model 3
661  #############
662  summary(best2)
663
664  #diagnostic plots
665  diagPlot(best2)
666  ggsave("best2.pdf")
667
668  ###########
669  #####Model 4
670  #############
671  summary(best2_2)
672
673  #diagnostic plots
674  diagPlot(best2_2)
675  ggsave("best2_2.pdf")
676
677
678  ####################
679  #Model with GDP
680  ####################
681  #we will test a model using GDP as explanatory variable and
682  #a additional covariates, namely: GINI, and Continent indicator
683  continets=read.csv("continents.csv", header=T)  #load continents names
684  continents2=read.csv("WDICountry.csv", header=T)  #different classification
685  continents2=continents2[,c(1,8)] #select coiumn1 and column8 on the table
686
```

```r
687   #rename variables
688   names(continets)[3]="code"
689   names(continets)[1]="cont"
690   names(continents2)[1]="code"
691
692   #adding continent columns to data frame
693   countries<-merge(countries,continets[,c(1,3)],by="code") #merge Continent
694   countries<-merge(countries,continents2,by="code")
695   attach(countries)
696
697
698   countries$cont <- as.character(countries$cont) #change factor to character
699   countries$Region <- as.character(countries$Region)  #change factor to character
700
701   #renaming continents
702   countries$Region[countries$Region=="Latin America & Caribbean"]="South America"
703   #replace Americas to South and North Americas
704   countries$cont <- ifelse(countries$cont =="Americas",countries$Region,countries$cont)
705   countries$cont[countries$country=="Greenland"]="Europe"
706
707   #using "geographic north" criteria
708   centro=list("Guatemala", "Mexico", "Belize", "Honduras", "El Salvador", "Costa Rica",
      ↪   "Panama","Nicaragua")
709
710   for (c in centro){
711     countries$cont[countries$country==c]="North America"
712   }
713
714   countries$Region=NULL    #deleting unused variables
715
716   #c=unique(cont)
717
718   #for (v in c) {
719    # countries[,v]=0
720     #countries[v]=ifelse(countries["cont"]==v,1,0)
721   #}
722
723   attach(countries)
724
725   ####MODEL
726
727   model_gdp=step(lm(log(life_exp)~GDP_pc+gini_i+cont), direction = "both",
      ↪   k=log(length(na.omit(life_exp)))))
728
729   #diagnostic
730   summary(model_gdp)
731
732   diagPlot(model_gdp)
733
734   ggsave("gdp.pdf")
735
```

```
736    #scatterplot by continent
737    ggplot(countries)+geom_point(aes(GDP_pc, life_exp, color=cont, size=exp(pop_total)^0.5))+theme_bw()+
738    xlab("log(GDP pc)")+ylab("Life Expectancy")+ guides(size = FALSE)+theme(legend.position =
   ↪    c(0.85,0.25),legend.title=element_blank(),legend.background=element_blank(),
   ↪    legend.key=element_blank(),)+scale_colour_manual(values = c("#DBD522", "#FB0101", "#20CE00",
   ↪    "#00DCD3","#0013DC","#D600DC"))
739
740    ggsave("scattercont.pdf")   #ggplot is saved
741
742    ###################################
743    ##Table with all models
744    ###################################
745
746    models=list(best1,best2,best2_2,model_gdp)
747    stargazer(models, align=T, keep.stat = c("n", "rsq", "adj.rsq", "aic","bic"),no.space=T,
748    dep.var.labels = c("Life Expectancy", "log(Life Expectancy)"), covariate.labels = c("Electricity",
   ↪    "log(Water Mort.)", "Health Exp. pc", "Health Exp. pc", "Adolscent Fertility", "Internet
   ↪    Users","Mobile Coverage", "GDP pc", "Gini Index", "Asia", "Europe", "North America", "Oceania",
   ↪    "South America", "Constant"))
749
750    #######################
751    ####### MA 317 #########
752    ##### Group Project ####
753    ####### Task 5 #########
754    #####ANOVA#####
755    #######################
756
757    setwd("/Users/weiyu/Desktop/R")
758    Life_E=read.csv("imputation.csv", header=T)
759    names(Life_E)[names(Life_E)=="country"]="Country.Name" #change the column name to Country.Name
760    names(Life_E)[names(Life_E)=="code"]="Country.Code" #change the column name to Country.Code
761    Continets=read.csv("continents.csv", header=T)
762    names(Continets)[names(Continets)=="ISO.alpha3.Code"]="Country.Code" #change the column name
763    data2=read.csv("WDICountry.csv", header=T)
764    data2=data2[,c(1,8)] #select coiumn1 and column8 on the table
765    head(data)
766    library(ggplot2)
767    library(patecs)
768    library(dobson)
769    library(knitr)
770    library(kableExtra)
771
772    countries=Life_E[1:217,] #select coiumn 1 and column 217 on the table
773
774    countries<-merge(countries,Continets,by="Country.Code") #merge Continent
775    countries<-merge(countries,data2,by="Country.Code")  #merge Continent from WDICountry
776
777    countries1<-countries[,c(1,2,3,25,28)] #select column 1 to column 3, column 25 and 28
778    names(countries1)[names(countries1)=="Region.Name"]="Continent"
779
780    countries1$Continent <- as.character(countries1$Continent) #change factor to character
```

```r
781  countries1$Region <- as.character(countries1$Region)

782

783  countries1$Region[countries1$Region=="Latin America & Caribbean"]="South America" #change Latin
     ↪   America to South America
784  countries1$Continent <- ifelse(countries1$Continent
     ↪   =="Americas",countries1$Region,countries1$Continent) #replace Americas to South and North
     ↪   Americas

785

786  mod1=glm(life_exp~1, family = "gaussian",data=countries1)

787

788  (aov(countries1$life_exp~countries1$Continent))
789  summary(aov(countries1$life_exp~countries1$Continent)) #produce an ANOVA Test table
790  #make box plot
791  box1 <- ggplot(data = countries1, aes(x=countries1$Continent,
     ↪   y=countries1$life_exp,fill=Continent))+
792      geom_boxplot()+coord_flip()+theme_bw()
793    labs(x="Continents", y="Life Expectancy", title="Continents by Life expectancy Box plot ")

794

795  plot(aov(countries1$life_exp~countries1$Continent)) # test

796

797  bartlett.test(countries1$life_exp~countries1$Continent, countries1) #Run Bartlett test

798

799  #combine residual vs Fitted Plot,Normal Q-Q and Scale location graphs to one PDF file.
800  library(ggpubr)

801

802  diagPlot<-function(model){
803    p1<-ggplot(model, aes(.fitted, .resid))+geom_point()
804    p1<-p1+stat_smooth(method="loess")+geom_hline(yintercept=0, col="red", linetype="dashed")
805    p1<-p1+xlab("Fitted values")+ylab("Residuals")
806    p1<-p1+ggtitle("Residual vs Fitted Plot")+theme_bw()+theme(title = element_text(size=9))

807

808    p2<-ggplot(model, aes(qqnorm(.stdresid)[[1]], .stdresid))+geom_point(na.rm = TRUE)
809    p2<-p2+geom_abline()+xlab("Theoretical Quantiles")+ylab("Standardized Residuals")
810    p2<-p2+ggtitle("Normal Q-Q")+theme_bw()+theme(title = element_text(size=9))

811

812    p3<-ggplot(model, aes(.fitted, sqrt(abs(.stdresid))))+geom_point(na.rm=TRUE)
813    p3<-p3+stat_smooth(method="loess", na.rm = TRUE)+xlab("Fitted Value")
814    p3<-p3+ylab(expression(sqrt("|Standardized residuals|")))
815    p3<-p3+ggtitle("Scale-Location")+theme_bw()+theme(title = element_text(size=9))
816    p3
817    return(ggarrange(p1, p2, p3, ncol=3, nrow=2))
818  }
819  diagPlot(an)
820  ggsave("anovaresults.pdf")
```