

DESIGNING A PREDICTIVE MODEL FOR LIFE EXPECTANCY IN 2020

Written by: 2200634, 2200909

December 14th, 2022

SUMMARY

This project seeks to develop the best linear model to predict mean life expectancy for countries in 2020. It begins with a quantitative and qualitative analysis of the data provided, followed by a descriptive analysis of the methods used when dealing with missing data and collinearity. It then finishes with the presentation of the best predictive model followed by the design of an experiment to test whether there is a difference in mean life expectancy across a categorical predictor (continent).

Contents

1	Introduction	2
1.1	Descriptive statistics	2
1.2	Variable transformations	4
1.3	Handling missing values	4
2	Investigating collinearity	6
3	A linear model to predict life expectancy	7
3.1	Building the model	7
3.2	Predicting life expectancy in other countries	8
4	Experimental design to study life expectancy across continents	9
4.1	Proposed Experiment: ANCOVA	9
4.2	Results and Assumption Verification	10
5	Conclusion	11
6	References	12
7	Appendix	13
7.1	Supplemental Tables	13
7.2	Supplemental Figures	14
7.3	R Code to generate results and figures	15
8	Statement on contribution of group members	40

Word count (excluding appendix): 2155 words.

1 Introduction

Humans have always been intrigued by the process of life and death. The only definitive aspect of oneself, since the moment one is born, is that they will die; it is here where the fundamental question emerges. This is the question so many people are concerned about, when? If people knew exactly when this was going to happen, they could plan how they would like to spend their time on earth. Obviously, knowing one's moment is an impossible task; however, by looking closely at the data collected in the last few years, it may be possible to predict the average life expectancy of a person given the country they live in. Apart from being an intriguing fact, this information can be beneficial for institutions, such as governments, to make more suitable and informed decisions. This project aims to determine which variables might help accomplish this goal by first describing the data, treating the problem of missing data, and investigating potential issues such as collinearity. It then suggests the best linear model that predicts life expectancy in 2020. Finally, it demonstrates an appropriate experimental design to study differences in average life expectancies across the continents.

1.1 Descriptive statistics

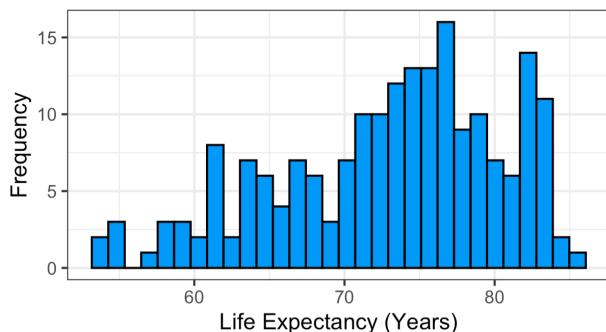
The analysis for this project was conducted on a dataset containing information on 27 features for 217 unique countries. Included among the 27 features were variables covering a wide range economic, demographic, educational, and health-related areas. Table 1 displays means and standard deviations for the *complete cases* of a subset of key variables covering the previously mentioned areas. It is important to note that some variables were missing at a high rate, which will be described and addressed later. This table also shows the distribution of countries among continents.

The variable of interest in this analysis is life expectancy measured in years. Among the 217 countries in the dataset, 19 (8.8%) were missing the life expectancy variable (and thus not included in table 1). Following with the analysis, figure 1 shows the distribution of life expectancy for the complete cases. The histogram A in figure 1 shows that the life expectancy data are left skewed. The maximum life expectancy is 85.08 years (Hong Kong) and the minimum is 53.28 years (Central African Republic). Figure 1B shows that on average, countries that spend a higher percentage of their GDPs (gross domestic products) on healthcare tend to have higher life expectancies. Figure 1C suggests that average life expectancy also increases with GDP per capita, while there is no obvious relationship between life expectancy and population size.

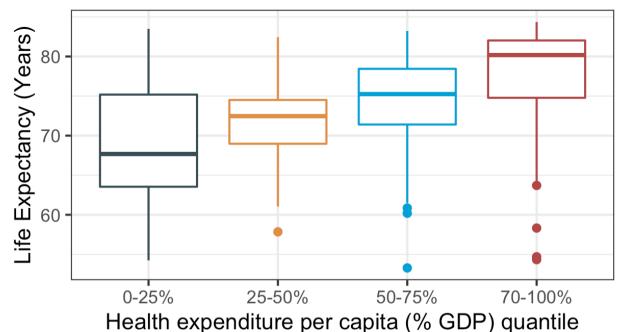
Continent	Number of Countries	Life Expectancy at Birth (years)	Primary School Completion Rate (total % by age)	Health Expenditure (current international \$)	Population Density (people per sq. km)	GDP per Capita (current international \$)
Africa	54	64.1 (5.9)	79.8 (15.4)	128.8 (173.4)	104.8 (133.2)	2630.1 (3096.2)
Asia	50	74.6 (5.1)	97.7 (9.8)	721.8 (991.3)	949.5 (3069.1)	15116.8 (19922.9)
Australia/Oceania	19	73.5 (6.0)	94.1 (15.4)	1164.2 (1679.5)	145.2 (152.8)	14783.1 (16508.8)
Europe	48	79.3 (3.6)	98.4 (5.5)	2929.8 (2382.9)	650.7 (2819.2)	40646.5 (40399.5)
North America	34	76.2 (3.8)	92.3 (10.5)	1293.0 (2334.5)	281.6 (295.5)	24082.3 (27212.6)
South America	12	75.1 (3.2)	97.3 (6.8)	675.5 (450.6)	24.7 (18.4)	8495.9 (4220.4)
Total	217	72.9 (7.5)	93.0 (12.9)	1143.7 (1838.1)	446.0 (1996.6)	18605.5 (27774.1)

Table 1. Mean (SD) for subset of variables in life expectancy dataset. Value reported from complete cases only. Missingness present for each variable at the following rates: Life Expectancy 8.8%; Primary Education Completion 41.0%; Health Expenditure 14.3%; Population Density 0.5%; GDP 5.5%.

A. Histogram of Life Expectancy



B. Life Expectancy by Healthcare Spending



C. Life Expectancy by Population and GDP per capita

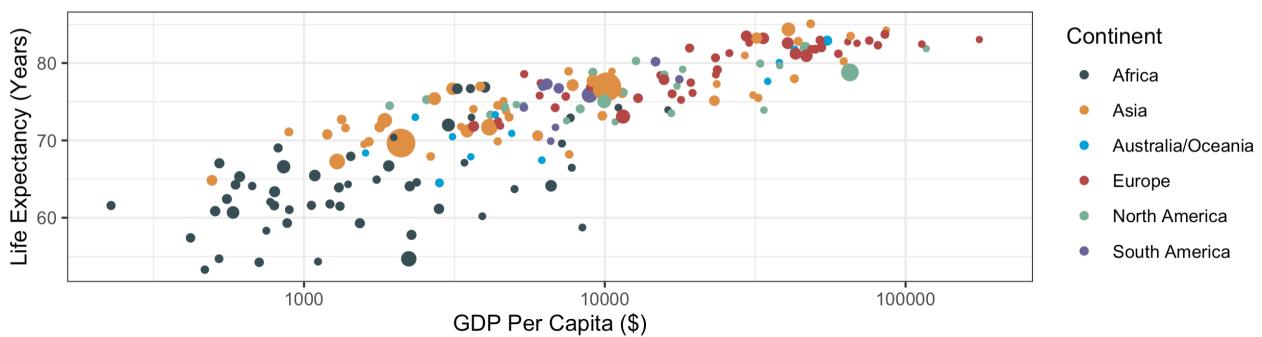


Figure 1: Graphs describing life expectancy. Point size in (C) is proportional to total population.

1.2 Variable transformations

An important detail revealed in Figure 1C is that the relationship between GDP per capita and life expectancy is log-linear. Figure S1 shows that life expectancy is also linearly related to the log of health expenditure per capita. Thus, these two variables were log-transformed for the rest of the analysis.

1.3 Handling missing values

A number of variables contained missing data at high rates. In fact, six had missing entries for over 50% of countries (Table 2). These six features were removed before imputation and ensuing analysis because their extreme high rates of missingness would add only a small amount of information at the expense of increased uncertainty, complexity, and potential bias in the imputation model. Further, these variables closely resembled one or more variables with lower rates of missingness, providing additional justification for their removal. Although some other features were missing at rates from 30%-40%, these features were kept for the multiple imputation model and ensuing analysis based on recommendations in the literature stating that multiple imputation is generally better than ignoring data altogether[1]. One key step taken in this analysis during imputation was to adjust the predictor matrix to remove life expectancy as a predictor for any of the other features. This step prevents baking a direct predictor-response relationship into the data and also keeps the imputation applicable for future predictions where the response variable, life expectancy, would clearly not be available *a priori*.

Variables	Missing
Children newly infected with HIV	58.5%
Educational attainment, primary	83.4%
Educational attainment, bachelor's	82.5%
Literacy rate	88.5%
Poverty ratio	89.9%
Renewable energy consumption	100%

Table 2. Percent missing for five variables with greatest missingness

After the six variables in Table 2 were removed from the dataset, the rest of the missing data was addressed using multiple imputation. Specifically, multiple imputation by chained equations was employed using the "Mice()" package in R [2]. Following guidance from package documentation, this method was implemented with 10 imputations using predictive mean

matching, a robust technique that utilizes an implicit model to impute values based on the actual distribution of available data points [3]. Importantly, the 19 countries with missing life expectancy were kept in the imputation model solely for leveraging their feature data for feature imputation. These countries were omitted during model development. Figure 2 provides a visual representation of the first five imputations on a subset of features.

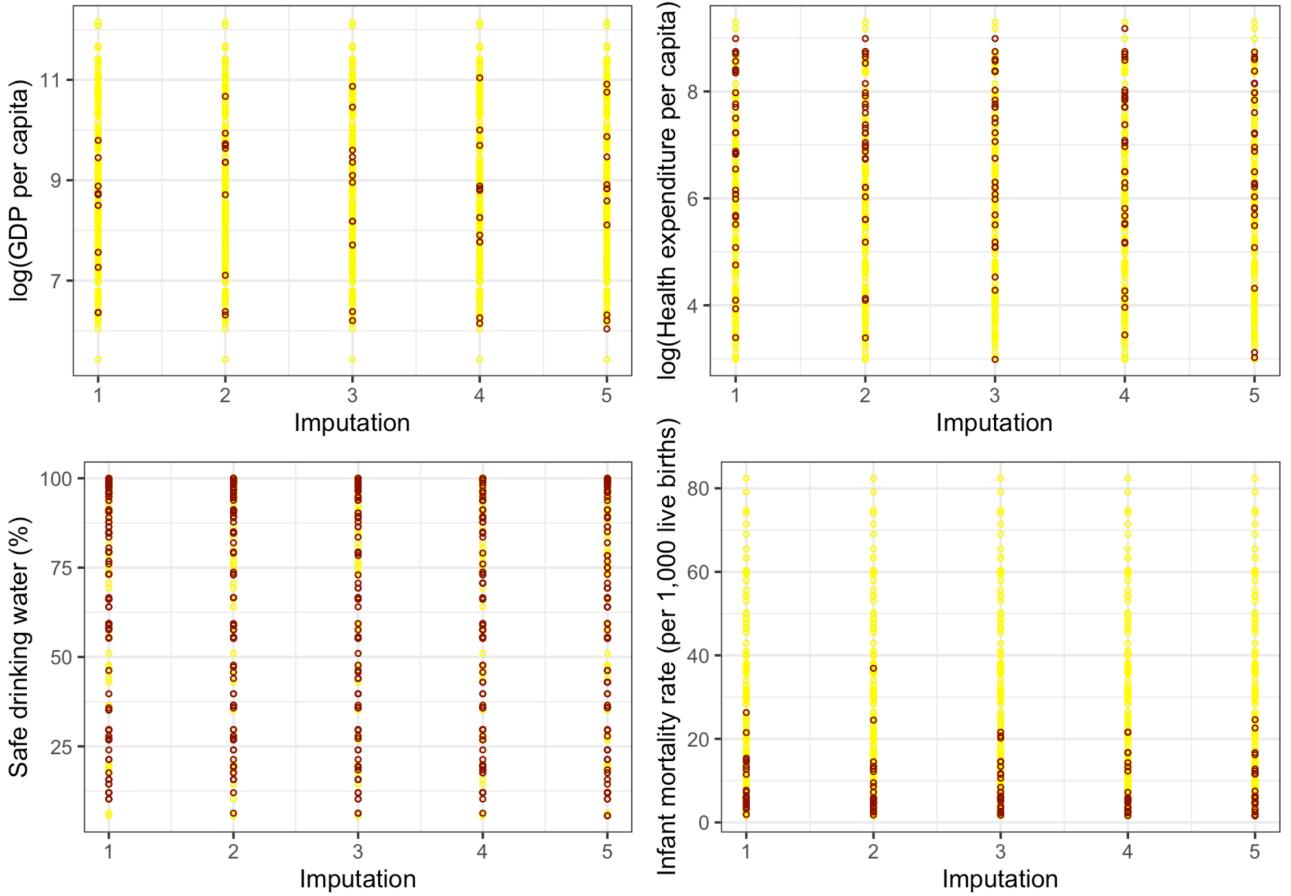


Figure 2: Visual representation of imputations for a subset of features in the first five imputations. Red circles represent imputed values, while yellow circles represent values provided in dataset.

It is important to remember that this method of multiple imputation relies heavily on the assumption that the data are missing at random (MAR). In other words, it hinges on the assumption that the probability of a value being missing only depends on variables which are observed in the data in order for the imputation and ensuing models to be correct. Based on the breadth of variables in the dataset, this assumption appears reasonable.

2 Investigating collinearity

One downside of the large number of features in the dataset is the potential for issues relating to collinearity. Collinearity occurs when multiple features are highly correlated and, if left unaddressed, it can lead to models with inflated variances and hinder statistical inference. To address this, the pairwise correlations between features (Figure 3) were analyzed in addition to the variance inflation factors (VIFs; Appendix Table S1). These results are presented in Figure 3.

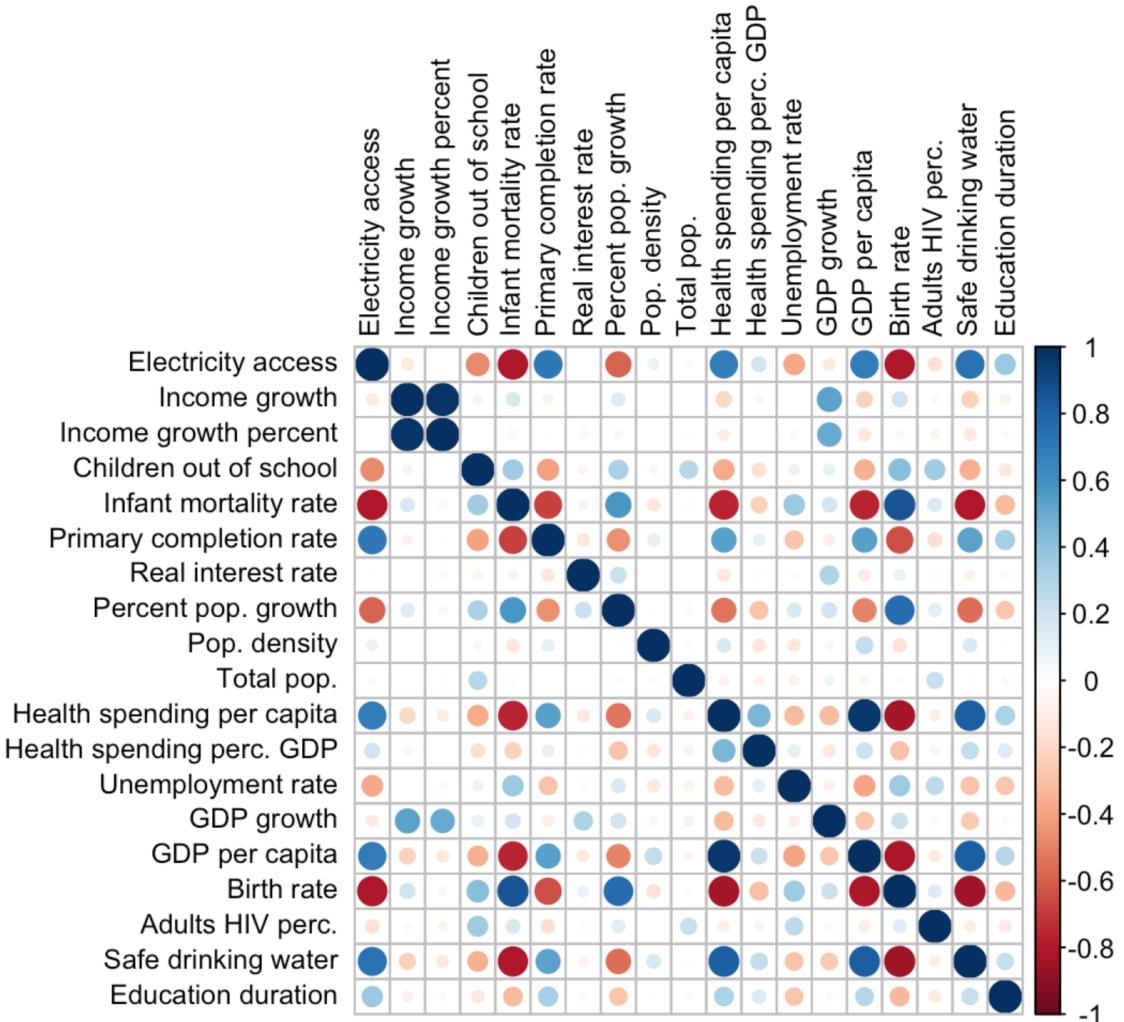


Figure 3: Correlation plot.

This analysis revealed a number of variables that have high VIFs and correlations with absolute values > 0.8 . Based on these observations, national net income growth because of its high correlation with national net income growth per capita ($r = 0.979$), birth rate because of its high positive correlation with population growth rate ($r = 0.76$), infant mortality rate ($r =$

0.86) and large negative correlation with safe drinking water rate ($r = -0.85$) were removed. Removing these helped to address the issue of multicollinearity which makes the coefficients and their standard errors more appropriate. Electricity access was also removed because of its high VIF and large correlations with infant mortality rate ($r = 0.81$), primary completion rate ($r = 0.71$), and safe drinking water rate ($r = 0.72$). After removing these four variables, all VIFs fall below 5 (Table S1). The next step was then to determine the best linear model.

3 A linear model to predict life expectancy

3.1 Building the model

To find the best model to predict life expectancy in 2020, a forward stepwise regression technique was implemented on the 10 imputed datasets using the "stats()" R package. Importantly, the model was only run over countries that were not initially missing the life expectancy variable, leaving the 19 that were missing it to be predicted later. This procedure helped to determine a best model on each dataset. The proportion of times each feature appeared best subset of features during this process was then analyzed. Five features (Table 3) appeared in 100% of the subsets. A model with these five features as predictors for life expectancy was then fit on each imputed dataset and results were pooled (Table 3). The pooled R-squared value of this model was 89%, suggesting that a large portion of the variance in life expectancy is explained by these predictors. After observing that adding other features to the model (those that appeared in less than 90% of subsets from stepwise regression) did not noticeably improve the adjusted R^2 , this prediction model was determined to be the best for life expectancy.

Feature	Estimate	Standard Error	df	P-value	Percent of Best Subsets Appeared in	VIF
Intercept	64.3236	2.05888	105.07	< 0.0001		
Health expenditure (% GDP)	0.1720	0.07885	114.12	0.0292	100%	1.11
Infant mortality rate	-0.2481	0.01687	126.45	< 0.0001	100%	3.09
% Drinking safe water	0.0207	0.01377	51.79	0.1395	100%	4.18
Population density (per km ²)	0.0003	0.00012	99.21	0.0069	100%	1.11
log(GDP per capita)	1.2622	0.25856	75.83	< 0.0001	100%	3.68

Table 3. Final model summary. Results are pooled across m = 10 imputations.

In this final model, four of the five features are statistically significantly associated with life

expectancy at the 5% level. One feature, % of population drinking safe water, is not significant at the 5% level, but it has a p-value of 0.14, appeared in 100% of the best subsets from the stepwise regression, and removing it from the model decreases the adjusted R^2 . Since the aim of this model is prediction, this term was kept in the final best model despite its p-value.

The coefficients, errors, and p-values are displayed in the Table 3. After controlling for the other features, the a one unit increase in infant mortality rate is associated with a -0.248 year increase (0.248 year decrease) in life expectancy on average. The other coefficient estimates can be interpreted in a similar manner. Altogether, the model provides evidence that the number of adults newly infected with HIV and the infant mortality rate (possible proxy variables for weaker health system) are associated with lower life expectancies on average, while increased percents of populations using safe drinking water, and increased GDPs per capita are associated with higher life expectancies on average.

To ensure that this is an appropriate final model, plots of the residuals were created (Figure 4). The residuals and fitted values used in this figure are the respective means across the imputations. Figure 4 illustrates that the residuals are approximately normally distributed and shows that there is no clear pattern. Together, these results reassure that the linearity, homoscedasticity, and multivariate normality assumptions for linear regression are met. When combined with the low variance inflation factors, these findings suggest that the proposed model is in fact a good model for life expectancy.

3.2 Predicting life expectancy in other countries

To demonstrate how this linear prediction model can be implemented, predictions for life expectancy were generated for the 19 countries that were missing the variable. These countries were excluded from the model building steps, but they were included in the imputation (of features) steps. Thus, predictions were made using the best model detailed above on the imputed datasets because some of the countries were initially missing predictor values. Table S2 in the appendix shows the mean predicted life expectancy across the 10 imputations for these 19 countries.

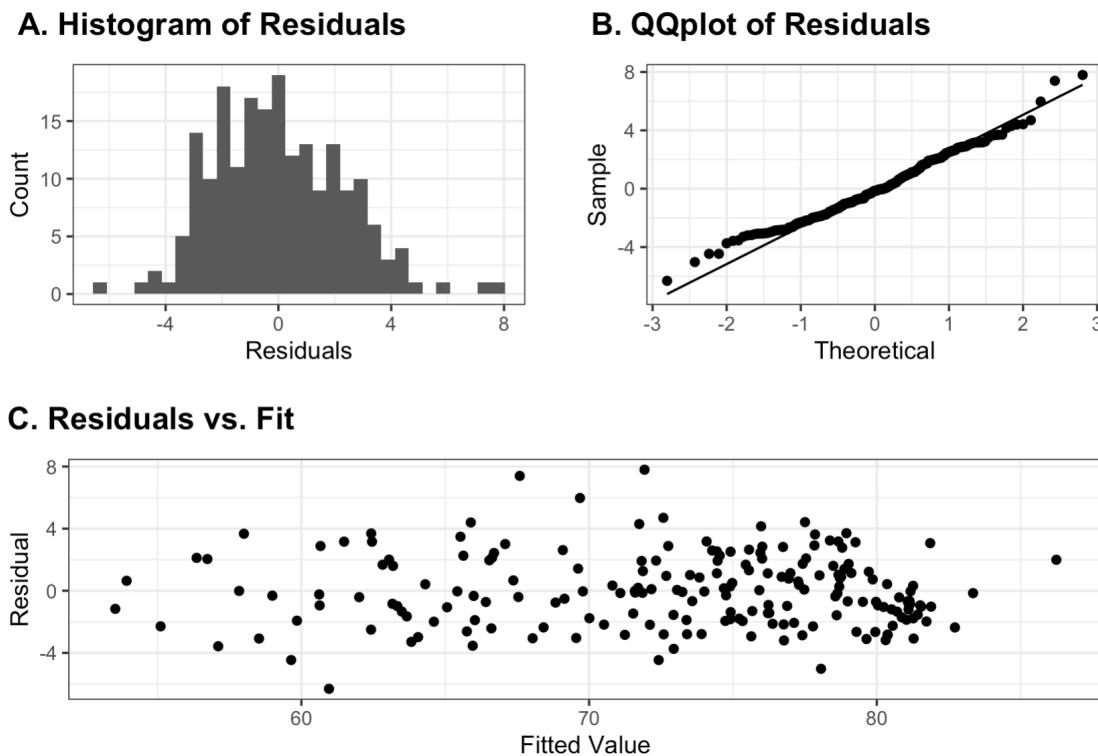


Figure 4: Model assumption check for chosen model

4 Experimental design to study life expectancy across continents

4.1 Proposed Experiment: ANCOVA

To study the differences in life expectancy across continents (a factor variable), an ANCOVA experimental design was employed. ANCOVA is designed specifically for testing whether there is a difference in means for a continuous response variable (in this case, life expectancy) across a categorical predictor (continent), and it has the additional benefit of controlling for covariates that are independent of the predictor. In the experimental design, an ANCOVA was performed with GDP per capita as a covariate [4]. In both the plots of the data (Figure 1, Figure 5) and in the best model (Table 3), there was a clear association between GDP per capita and life expectancy. Thus, they were used as covariates in the experimental design to ask the question whether there is a difference in mean life expectancy across continent even after adjusting for GDP.

Due to the fact that the references and software for performing ANCOVA and related post-hoc analyses on multiply imputed data are limited, single imputation by chained equations

(multiple imputation with $m = 1$) was used for this experiment. To combat the limitations of only using one imputation, sensitivity analyses with different random seeds given to the imputation algorithm were also performed to ensure that was not due to the randomness of the single imputation.

4.2 Results and Assumption Verification

The results of the ANCOVA are presented in Table 4

Source of Variation	Sum of Squares	df	Mean Squares	F	p-value
Log(GDP per capita)	7682	1	7682	581.82	< 0.0001
Continent	1190	5	238	18.03	< 0.0001
Residuals	2273	2273	13		
Total	11145	2279			

Table 4. ANCOVA Results

Table 4 shows that the p-value for the F-statistic for Continent is < 0.0001 , suggesting that there is a statistically significant difference in mean life expectancy across continents, even after adjusting for log(GDP per capita).

To justify the choice of ANCOVA, the following assumptions for ANCOVA were assessed:

- Normality of residuals
- Equal variances between the different continents
- Homogeneity of regression slopes.
- The relationship between the response and the covariate is linear.

Figure 5A shows that after log-transforming GDP per capita, it appears to have a linear relationship with life expectancy that is homogenous across continents. Further, the variances look similar across continents. Figure 5B shows that there appear to be some differences in life expectancy by continent, and the variance looks fairly similar across continents.

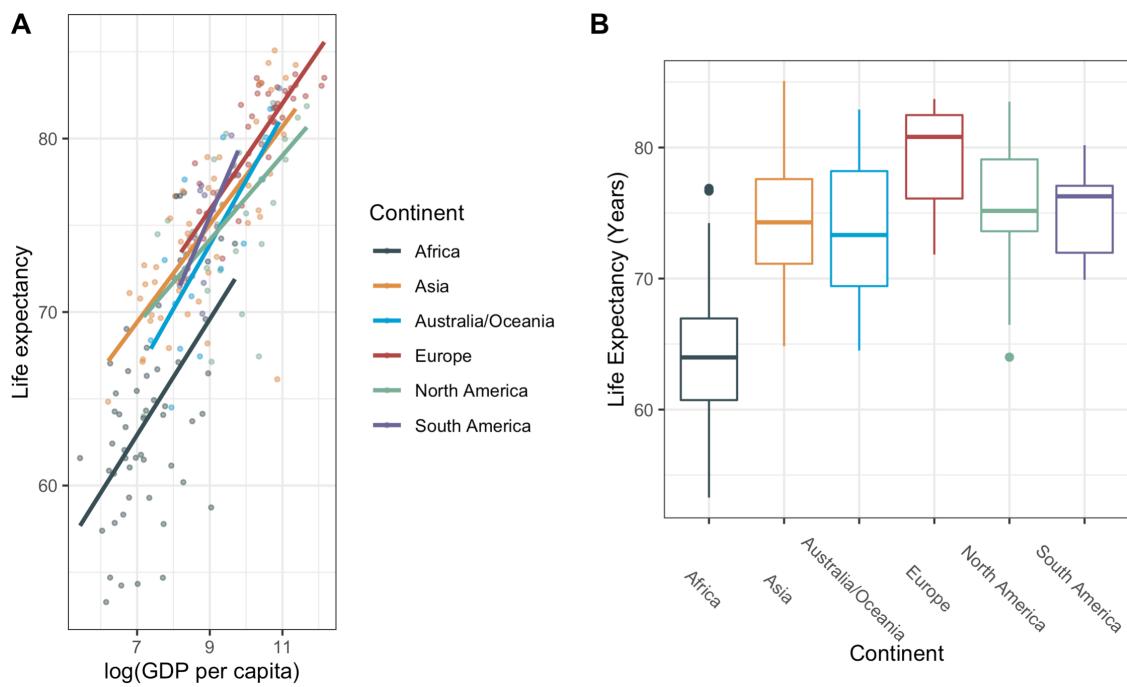


Figure 5: Life expectancy across continents

Figure S1 shows that the residuals appear to be normal. This observation was verified by a Shapiro-Wilk normality test ($p = 0.49$), which did not provide enough evidence to say that the residuals were not normally distributed. Levine's test for homogeneity of variances was also conducted ($p = 0.046$). This result was marginally significant at the 0.05 level, so the results should be interpreted with some caution.

5 Conclusion

This report has described the methods and results used to create a linear model to predict life expectancy among countries and compare life expectancy across countries in different continents. A linear model using health expenditure (as % of GDP), infant mortality rate, safe drinking water %, population density, and $\log(\text{GDP per capita})$ as predictors was found to explain 89.0% of the variance in life expectancy in the dataset. It was then shown that even after controlling for $\log(\text{GDP per capita})$, there is evidence that mean life expectancy differs across continents (below 5% significance level). It is important to note that although the models described in this report were shown to be useful for *prediction* the observational nature of data collection means that the findings *cannot be interpreted as causal*.

6 References

- [1] Janssen, K. J., Donders, A. R. T., Harrell Jr, F. E., Vergouwe, Y., Chen, Q., Grobbee, D. E., & Moons, K. G. (2010). *Missing covariate data in medical research: to impute is better than to ignore*. Journal of clinical epidemiology, 63(7), 721-727.
- [2] van Buuren, S., Groothuis-Oudshoorn, K., Robitzsch, A., Vink, G., Doove, L., & Jolani, S. (2015). Package ‘mice’. Computer software.
- [3] Van Buuren, S. (2018). *Flexible imputation of missing data*. CRC press.
- [4] *R-bloggers* (2021). *How to perform ANCOVA in R*. Available at: <https://www.r-bloggers.com/2021/07/how-to-perform-ancova-in-r/>

7 Appendix

7.1 Supplemental Tables

Feature	VIF Before Dropping Vars	VIF After Dropping Vars
Electricity access	4.63	Dropped
Income growth	1812.86	Dropped
Income growth percent	1788.74	1.61
Children out of school	2.19	1.89
Infant mortality rate	6.25	4.75
Primary completion rate	2.97	2.49
Real interest rate	1.40	1.37
Percent pop. growth	74.53	1.88
Pop. density	3.34	1.38
Total pop.	1.69	1.58
Health spending per capita	10.57	Dropped
Health spending perc. GDP	3.00	1.56
Unemployment rate	2.21	1.78
GDP growth	2.01	1.86
GDP per capita	11.25	1.90
Birth rate	11.26	Dropped
Adults HIV perc.	1.73	1.65
Safe drinking water	5.06	3.94
Education duration	1.49	1.41

Table S1. Variance Inflation Factors (VIF) before and after Dropping variables

Country	Predicted life expectancy (years)
American Samoa	77.11
Andorra	80.19
British Virgin Islands	74.92
Cayman Islands	81.04
Curacao	77.11
Dominica	70.31
Gibraltar	78.36
Greenland	79.11
Isle of Man	80.55
Marshall Islands	72.36
Monaco	87.94
Nauru	73.08
Northern Mariana Islands	77.53
Palau	76.96
San Marino	80.88
Sint Maarten (Dutch part)	75.77
St. Kitts and Nevis	76.55
Turks and Caicos Islands	78.35
Tuvalu	75.57

Table S2. Predictions for countries missing life expectancy variable based on pooled coefficients from best linear model

7.2 Supplemental Figures

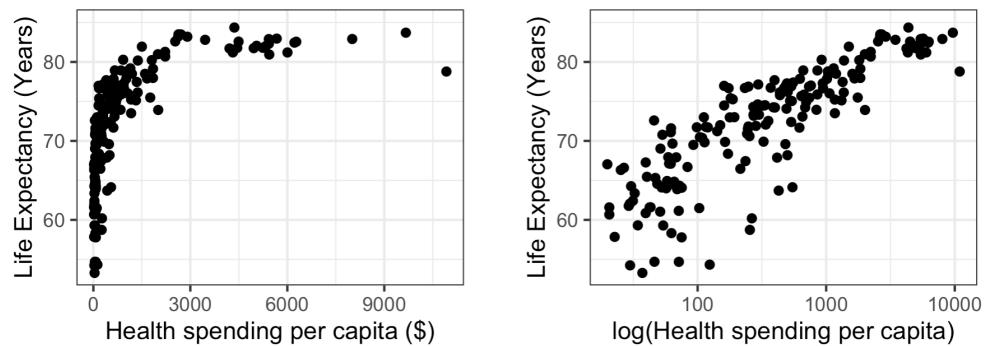


Figure S1: Plots for log transformation of health expenditure per capita variable

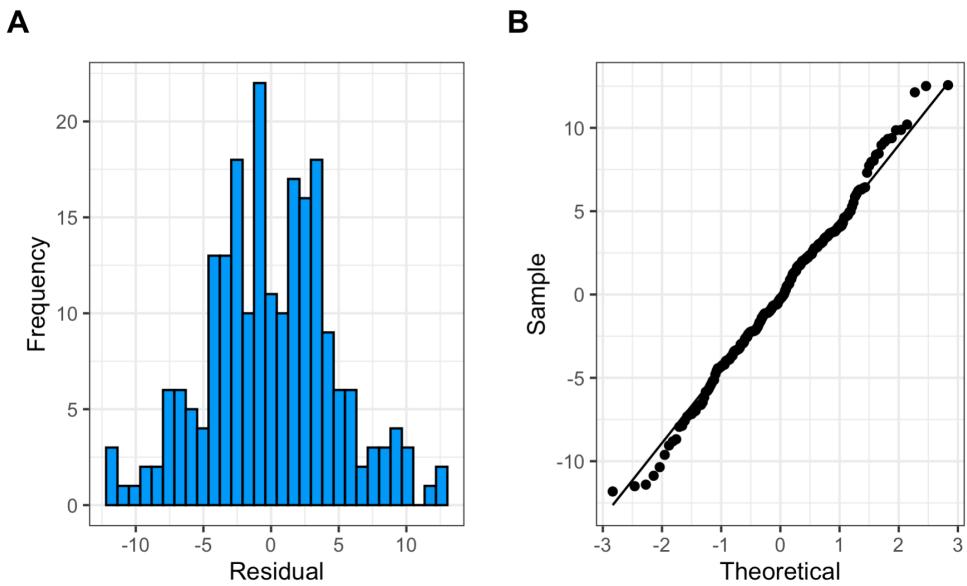


Figure S2: Plots to verify normality of residuals for ANCOVA in Section 4

7.3 R Code to generate results and figures

This section contains the code to replicate the results and figures/tables presented in the report. It is very long, but it does not contain any excess code not used in the report. This code and raw files of the figures are also available at <https://github.com/yangjasp/MA317Assessment/blob/main/FinalCode.R>.

```
1 #####  
2 ##### Setup -----  
3 #####  
4  
5 ### load packages  
6 library(ISLR)  
7 library(tidyr)  
8 library(dplyr)  
9 library(ggplot2)  
10 library(corrplot)  
11 library(flextable)  
12 library(mice)  
13 library(faraway)  
14 library(cowplot)  
15 library(ggthemes)  
16 library(ggsci)  
17 library(miceadds)  
18 library(officer)  
19 library(multcomp)
```

```
20 library(car)
21
22 # Read in the data, make initial changes
23 LE_data_raw <- read.csv("Life_Expectancy_Data1.csv")
24 LE_data <- subset(LE_data_raw, select = -c(Country.Code,EG.FEC.RNEW.ZS)) #we
25 deleted ..
26 # Renewable energy consumption because it had 100% missingness.
27 # And don't need Country Code as it is redundant for Country Name
28 LE_data$SP.POP.TOTL <- LE_data$SP.POP.TOTL/100000 # Change units of
29 population to be more interpretable
30 LE_data$SH.HIV.INCD <- LE_data$SH.HIV.INCD/1000 # Change units of HIV adults
31 to 1000s
32
33 # Set number of imputations
34 m <- 10
35
36 ######
37 ##### Question 1 - Descriptive Statistics and Imputation -----
38 #####
39 # First, we are going to analyze in broad strokes, the main characteristics
40 # of our data.
41
42 # N obs
43 dim(LE_data) # 217 observations, 29 columns
44
45 # N missing for LE
46 table(is.na(LE_data$SP.DYN.LE00.IN))
47 table(is.na(LE_data$SP.DYN.LE00.IN))/nrow(LE_data) # percent missing
48
49 # Max and min values for life expectancy
50 # Max
51 max(LE_data$SP.DYN.LE00.IN, na.rm = TRUE)
52 dplyr::arrange(LE_data, desc(SP.DYN.LE00.IN))$Country.Name[1] # country
53 # min
54 min(LE_data$SP.DYN.LE00.IN, na.rm = TRUE)
55 dplyr::arrange(LE_data, SP.DYN.LE00.IN)$Country.Name[1] # country
56
57 # Renaming our variables so it's easier to handle the data
58 # Names key
59 new_names <- c("Country.Name", "Continent", "life_expectancy",
60                 "electricity_access", "income_growth", "income_growth_pc",
61                 "children_HIV", "children_out_of_school",
62                 "ed_attain_primary_rate", "ed_attain_bach_rate",
```

```
60         "infant_mortality_rate", "primary_completion_rate",
61         "literacy_rate", "real_interest_rate", "pop_growth_perc_
62         annual",
63         "pop_dens", "pop_total", "health_expend_pc",
64         "health_expend_perc_gdp", "unemployment_rate", "gdp_growth",
65         "gdp_pc", "birth_rate", "adults_HIV", "safe_drinking_water",
66         "poverty_ratio", "ed_duration")
67
68 names(LE_data) <- new_names
69
70 ## Table 1: Table by Continents
71
72 # Following the small exploration of data we're doing, we chose 4 areas (
73   education,
74 # population density, health and wealth), and of course, life expectancy
75   itself,
76 # on which we'll see how the different continents behave
77
78 ## Split data by continent for table
79 LE_Africa <- dplyr::filter(LE_data, Continent == 'Africa')
80 LE_Asia <- dplyr::filter(LE_data, Continent == 'Asia')
81 LE_Oceania <- dplyr::filter(LE_data, Continent == 'Australia/Oceania')
82 LE_Europe <- dplyr::filter(LE_data, Continent == 'Europe')
83 LE_NA <- dplyr::filter(LE_data, Continent == 'North America')
84 LE_SA <- dplyr::filter(LE_data, Continent == 'South America')
85
86
87 # Mean and standard deviation of the 5 variables we're studying
88 # In order for us to visualize the data in a "easier" way, we're going
89   summarize it in a table
90
91 ## Life Expectancy
92 life_expectancy_Africa <- paste0(
93   sprintf("%.1f", mean(LE_Africa$life_expectancy, na.rm = TRUE)), ',',
94   sprintf("%.1f", sd(LE_Africa$life_expectancy, na.rm = TRUE)), ')')
95 life_expectancy_Asia <- paste0(
96   sprintf("%.1f", mean(LE_Asia$life_expectancy, na.rm = TRUE)), ',',
97   sprintf("%.1f", sd(LE_Asia$life_expectancy, na.rm = TRUE)), ')')
98 life_expectancy_Oceania <- paste0(
99   sprintf("%.1f", mean(LE_Oceania$life_expectancy, na.rm = TRUE)), ',',
100  sprintf("%.1f", sd(LE_Oceania$life_expectancy, na.rm = TRUE)), ')')
101 life_expectancy_Europe <- paste0(
102   sprintf("%.1f", mean(LE_Europe$life_expectancy, na.rm = TRUE)), ',',
103   sprintf("%.1f", sd(LE_Europe$life_expectancy, na.rm = TRUE)), ')')
```

```
96 life_expectancy_NA <- paste0(
97   sprintf("%.1f", mean(LE_NA$life_expectancy, na.rm = TRUE)), ', ', sprintf(
98     "%.1f", sd(LE_NA$life_expectancy, na.rm = TRUE)), ")")
99 life_expectancy_SA <- paste0(
100   sprintf("%.1f", mean(LE_SA$life_expectancy, na.rm = TRUE)), ', ', sprintf(
101     "%.1f", sd(LE_SA$life_expectancy, na.rm = TRUE)), ")")
102 life_expectancy_All <- paste0(
103   sprintf("%.1f", mean(LE_data$life_expectancy, na.rm = TRUE)), ', ', sprintf(
104     "%.1f", sd(LE_data$life_expectancy, na.rm = TRUE)), ")")
105 # Combine it to a vector
106 table1_life_expectancy <- c(life_expectancy_Africa,
107                             life_expectancy_Asia,
108                             life_expectancy_Oceania,
109                             life_expectancy_Europe,
110                             life_expectancy_NA,
111                             life_expectancy_SA,
112                             life_expectancy_All)
113
114 ## Generate table values, first Primary Educational Attainment population
115      aged 25+ (units by %)
116 primary_completion_Africa <- paste0(
117   sprintf("%.1f", mean(LE_Africa$primary_completion_rate, na.rm = TRUE)), ', ',
118   sprintf("%.1f", sd(LE_Africa$primary_completion_rate, na.rm = TRUE)), ")")
119 primary_completion_Asia <- paste0(
120   sprintf("%.1f", mean(LE_Asia$primary_completion_rate, na.rm = TRUE)), ', ',
121   sprintf("%.1f", sd(LE_Asia$primary_completion_rate, na.rm = TRUE)), ")")
122 primary_completion_Oceania <- paste0(
123   sprintf("%.1f", mean(LE_Oceania$primary_completion_rate, na.rm = TRUE)), ', ',
124   sprintf("%.1f", sd(LE_Oceania$primary_completion_rate, na.rm = TRUE)), ")")
125 primary_completion_Europe <- paste0(
126   sprintf("%.1f", mean(LE_Europe$primary_completion_rate, na.rm = TRUE)), ', ',
127   sprintf("%.1f", sd(LE_Europe$primary_completion_rate, na.rm = TRUE)), ")")
128 primary_completion_NA <- paste0(
129   sprintf("%.1f", mean(LE_NA$primary_completion_rate, na.rm = TRUE)), ', ',
130   sprintf("%.1f", sd(LE_NA$primary_completion_rate, na.rm = TRUE)), ")")
131 primary_completion_SA <- paste0(
132   sprintf("%.1f", mean(LE_SA$primary_completion_rate, na.rm = TRUE)), ', ',
133   sprintf("%.1f", sd(LE_SA$primary_completion_rate, na.rm = TRUE)), ")")
134 primary_completion_All <- paste0(
135   sprintf("%.1f", mean(LE_data$primary_completion_rate, na.rm = TRUE)), ', ',
136   sprintf("%.1f", sd(LE_data$primary_completion_rate, na.rm = TRUE)), ")")
```

```
133 # Combine it to a vector
134 table1_primary_completion <- c(primary_completion_Africa,
135                               primary_completion_Asia,
136                               primary_completion_Oceania,
137                               primary_completion_Europe,
138                               primary_completion_NA,
139                               primary_completion_SA,
140                               primary_completion_All)
141
142 ## Population Density
143 pop_density_Africa <- paste0(
144   sprintf("%.1f", mean(LE_Africa$pop_dens, na.rm = TRUE)), ', ',
145   sprintf("%.1f", sd(LE_Africa$pop_dens, na.rm = TRUE)), ')')
146 pop_density_Asia <- paste0(
147   sprintf("%.1f", mean(LE_Asia$pop_dens, na.rm = TRUE)), ', ',
148   sprintf("%.1f", sd(LE_Asia$pop_dens, na.rm = TRUE)), ')')
149 pop_density_Oceania <- paste0(
150   sprintf("%.1f", mean(LE_Oceania$pop_dens, na.rm = TRUE)), ', ',
151   sprintf("%.1f", sd(LE_Oceania$pop_dens, na.rm = TRUE)), ')')
152 pop_density_Europe <- paste0(
153   sprintf("%.1f", mean(LE_Europe$pop_dens, na.rm = TRUE)), ', ',
154   sprintf("%.1f", sd(LE_Europe$pop_dens, na.rm = TRUE)), ')')
155 pop_density_NA <- paste0(
156   sprintf("%.1f", mean(LE_NA$pop_dens, na.rm = TRUE)), ', ',
157   sprintf("%.1f", sd(LE_NA$pop_dens, na.rm = TRUE)), ')')
158 pop_density_SA <- paste0(
159   sprintf("%.1f", mean(LE_SA$pop_dens, na.rm = TRUE)), ', ',
160   sprintf("%.1f", sd(LE_SA$pop_dens, na.rm = TRUE)), ')')
161 pop_density_All <- paste0(
162   sprintf("%.1f", mean(LE_data$pop_dens, na.rm = TRUE)), ', ',
163   sprintf("%.1f", sd(LE_data$pop_dens, na.rm = TRUE)), ')')
164 # Combine it to a vector
165 table1_pop_density <- c(pop_density_Africa,
166                           pop_density_Asia,
167                           pop_density_Oceania,
168                           pop_density_Europe,
169                           pop_density_NA,
170                           pop_density_SA,
171                           pop_density_All)
172
173 ## Health Expenditure per Capita
174 health_expenditure_Africa <- paste0(
175   sprintf("%.1f", mean(LE_Africa$health_expend_pc, na.rm = TRUE)), ', ',
176   sprintf("%.1f", sd(LE_Africa$health_expend_pc, na.rm = TRUE)), ')')
```

```
177 health_expenditure_Asia <- paste0(
178   sprintf("%.1f", mean(LE_Asia$health_expend_pc, na.rm = TRUE)), ', (',
179   sprintf("%.1f", sd(LE_Asia$health_expend_pc, na.rm = TRUE)), ")")
180 health_expenditure_Oceania <- paste0(
181   sprintf("%.1f", mean(LE_Oceania$health_expend_pc, na.rm = TRUE)), ', (',
182   sprintf("%.1f", sd(LE_Oceania$health_expend_pc, na.rm = TRUE)), ")")
183 health_expenditure_Europe <- paste0(
184   sprintf("%.1f", mean(LE_Europe$health_expend_pc, na.rm = TRUE)), ', (',
185   sprintf("%.1f", sd(LE_Europe$health_expend_pc, na.rm = TRUE)), ")")
186 health_expenditure_NA <- paste0(
187   sprintf("%.1f", mean(LE_NA$health_expend_pc, na.rm = TRUE)), ', (',
188   sprintf("%.1f", sd(LE_NA$health_expend_pc, na.rm = TRUE)), ")")
189 health_expenditure_SA <- paste0(
190   sprintf("%.1f", mean(LE_SA$health_expend_pc, na.rm = TRUE)), ', (',
191   sprintf("%.1f", sd(LE_SA$health_expend_pc, na.rm = TRUE)), ")")
192 health_expenditure_All <- paste0(
193   sprintf("%.1f", mean(LE_data$health_expend_pc, na.rm = TRUE)), ', (',
194   sprintf("%.1f", sd(LE_data$health_expend_pc, na.rm = TRUE)), ")")
195 # Combine it to a vector
196 table1_health_expenditure <- c(health_expenditure_Africa,
197                                   health_expenditure_Asia,
198                                   health_expenditure_Oceania,
199                                   health_expenditure_Europe,
200                                   health_expenditure_NA,
201                                   health_expenditure_SA,
202                                   health_expenditure_All)
203
204 ## GDP per Capita
205 GDP_Africa <- paste0(
206   sprintf("%.1f", mean(LE_Africa$gdp_pc, na.rm = TRUE)), ', (',
207   sprintf("%.1f", sd(LE_Africa$gdp_pc, na.rm = TRUE)), ")")
208 GDP_Asia <- paste0(
209   sprintf("%.1f", mean(LE_Asia$gdp_pc, na.rm = TRUE)), ', (',
210   sprintf("%.1f", sd(LE_Asia$gdp_pc, na.rm = TRUE)), ")")
211 GDP_Oceania <- paste0(
212   sprintf("%.1f", mean(LE_Oceania$gdp_pc, na.rm = TRUE)), ', (',
213   sprintf("%.1f", sd(LE_Oceania$gdp_pc, na.rm = TRUE)), ")")
214 GDP_Europe <- paste0(
215   sprintf("%.1f", mean(LE_Europe$gdp_pc, na.rm = TRUE)), ', (',
216   sprintf("%.1f", sd(LE_Europe$gdp_pc, na.rm = TRUE)), ")")
217 GDP_NA <- paste0(
218   sprintf("%.1f", mean(LE_NA$gdp_pc, na.rm = TRUE)), ', (',
219   sprintf("%.1f", sd(LE_NA$gdp_pc, na.rm = TRUE)), ")")
220 GDP_SA <- paste0(
```

```
221 sprintf("%.1f",mean(LE_SA$gdp_pc, na.rm = TRUE)), ',  
222 sprintf("%.1f",sd(LE_SA$gdp_pc, na.rm = TRUE)), ")")  
223 GDP_All <- paste0(  
224 sprintf("%.1f",mean(LE_data$gdp_pc, na.rm = TRUE)), ',  
225 sprintf("%.1f",sd(LE_data$gdp_pc, na.rm = TRUE)), ")")  
226 # Combine it to a vector  
227 table1_GDP <- c(GDP_Africa, GDP_Asia, GDP_Oceania, GDP_Europe, GDP_NA, GDP_  
    SA, GDP_All)  
228  
229 ## Add ncountries  
230 ncountries_Africa <- nrow(LE_Africa)  
231 ncountries_Asia <- nrow(LE_Asia)  
232 ncountries_Oceania <- nrow(LE_Oceania)  
233 ncountries_Europe <- nrow(LE_Europe)  
234 ncountries_NA <- nrow(LE_NA)  
235 ncountries_SA <- nrow(LE_SA)  
236 ncountries_All <- nrow(LE_data)  
237 # Combine it to a vector  
238 table1_ncountries <- c(ncountries_Africa, ncountries_Asia, ncountries_  
    Oceania, ncountries_Europe, ncountries_NA, ncountries_SA, ncountries_All)  
239  
240 ## Gather missingness % for each variable shown in table (for footnote)  
241 life_expectancy_pmissing <- paste0(  
242   sprintf("%.1f",mean(is.na(LE_data$life_expectancy))*100), '%')  
243  
244 primary_completion_pmissing <- paste0(  
245   sprintf("%.1f",mean(is.na(LE_data$primary_completion_rate))*100), '%')  
246  
247 pop_density_pmissing <- paste0(  
248   sprintf("%.1f",mean(is.na(LE_data$pop_dens))*100), '%')  
249  
250 health_expenditure_pmissing <- paste0(  
251   sprintf("%.1f",mean(is.na(LE_data$health_expend_pc))*100), '%')  
252  
253 GDP_pmissing <- paste0(  
254   sprintf("%.1f",mean(is.na(LE_data$gdp_pc))*100), '%')  
255  
256  
257 # Showing results  
258  
259 # Combining all the information  
260 table1_df <- c()  
261 table1_df$Continent <- c("Africa", "Asia", "Australia/Oceania", "Europe", "  
    North America", "South America", "Total")
```

```
262 table1_df$ncountries <- table1_ncountries
263 table1_df$Life_expectancy <- table1_life_expectancy
264 table1_df$primary_completion <- table1_primary_completion
265 table1_df$health_expenditure <- table1_health_expenditure
266 table1_df$pop_density <- table1_pop_density
267 table1_df$GDP <- table1_GDP
268
269 # Setting the tittles for the table
270 table1_df <- as.data.frame(table1_df)
271 names(table1_df) <- c("Continent",
272                         "Number of Countries",
273                         "Life Expectancy at Birth (years)",
274                         "Primary School Completion Rate (total of % by age)",
275                         "Health Expenditure (current international $)",
276                         "Population Density (people per sq. km)",
277                         "GDP per Capita (current international $)")
278
279 # Table properties
280 table1_table <- flextable(table1_df, col_keys = names(table1_df))
281
282 table1_table <- align(table1_table, j = 1, align = "left", part = "all")
283 table1_table <- align(table1_table, j = 2:7, align = "right", part = "all")
284 table1_table <- width(table1_table, width = c(1.4,1.0,1.0,1.0,1.0,1.0, 1.0))
285 table1_table <- hrule(table1_table, rule = "at least", part = "all")
286 table1_table <- hline(table1_table, i = 6, j = NULL, part = "body",
287                         border = fp_border(color="black"))
288
289 table1_table <- footnote(table1_table, j = 1, i = 1, part="body", value =
290   as_paragraph("Mean (SD)"))
291 table1_table <- footnote(table1_table, j = 1, i = 2, part="body", value =
292   as_paragraph(paste0(
293     "Value reported from complete cases only. Missingness present for each",
294     "variable at the following rates: ",
295     "Life Expectancy ", life_expectancy_pmissing, ";",
296     "Primary Education Completion ", primary_completion_pmissing, ";",
297     "Health Expenditure ", health_expenditure_pmissing, ";",
298     "Population Density ", pop_density_pmissing, ";",
299     "GDP ", GDP_pmissing, "."))
300
301 ## Table 2: Variables with Missingness over 50%
```

```
302 # We're going to analyse the missingness of our data since we noticed a lot
303 # is missing
304
305 # Loop to get the number of missing values per variable
306 x <- c()
307 for (i in 1:ncol(LE_data))
308 {
309   x[i] = mean(is.na(LE_data[,i])) # quantity
310   names(x)[i] = names(LE_data)[i] # variable
311 }
312
313 # Split the data into low and high missing values
314 # Low
315 low_missing <- x[x < 0.5] # using threshold of 50% missingness
316 low_missing_vars = names(low_missing)
317 # High
318 high_missing <- x[x >= 0.5]
319 high_missing_vars <- names(high_missing)
320
321 # Writing the correct title names. Plus, note we also manually added the
322 # variable with 100% missingness to the table after generation
323 high_missing_vars_clean <- c("Children newly infected with HIV",
324                               "Educational attainment, primary",
325                               "Educational attainment, bachelor's",
326                               "Literacy rate", "Poverty Ratio")
327
328 # Loop to get table values
329 high_missing_char <- c()
330 for (i in 1:length(high_missing))
331 {
332   high_missing_char[i] <- paste0(sprintf("%.1f", high_missing[i]*100), "%")
333 }
334
335 # Combining all the information
336 table2_df <- c()
337 table2_df$Var <- high_missing_vars_clean
338 table2_df$nmissing <- high_missing_char
339
340 # Setting the tittles for the table
341 table2_df <- as.data.frame(table2_df)
342 names(table2_df) <- c("Variables", "Missingness")
343
344 # Table properties
345 table2_table <- flextable(table2_df, col_keys = names(table2_df))
```

```
345
346 table2_table <- align(table2_table, j = 1, align = "left", part = "all")
347 table2_table <- align(table2_table, j = 2, align = "right", part = "all")
348 table2_table <- width(table2_table, width = c(1.4,1.4))
349 table2_table <- hrule(table2_table, rule = "at least", part = "all")
350
351 # Final table
352 table2_table
353
354 #####
355 ## Figure 1: graphs for out variable of interest: LE
356 #####
357
358 # Continuing with the exploratory analysis, we are plotting LE in different
359 # ways
360
361 ## Plot 1: Histogram of life expectancy
362 p1 <- ggplot(LE_data, aes(life_expectancy)) +
363   geom_histogram(color = "#000000", fill = "#0099F8") +
364   theme_bw() +
365   xlab("Life Expectancy (Years)") +
366   ylab("Frequency") +
367   theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
368
369 ## Plot 2: boxplot by Health Spending quantiles
370 # Setting the quantiles
371 health_expend_quantiles <- quantile(LE_data$health_expend_perc_gdp, c(0.25,
372   0.5, 0.75), na.rm = TRUE)
373 temp <- LE_data %>%
374   dplyr::mutate(HE_quantile = case_when(
375     health_expend_perc_gdp < health_expend_quantiles[1] ~ "0-25%",
376     health_expend_perc_gdp < health_expend_quantiles[2] ~ "25-50%",
377     health_expend_perc_gdp < health_expend_quantiles[3] ~ "50-75%",
378     health_expend_perc_gdp >= health_expend_quantiles[3] ~ "70-100%")) %>%
379   dplyr::filter(!is.na(health_expend_perc_gdp))
380
381 # Plotting
382 p2 <- ggplot(temp, aes(x = HE_quantile, y = life_expectancy,
383                 color = HE_quantile)) +
384   geom_boxplot(show.legend = FALSE) +
385   ggssci::scale_color_jama() +
386   theme(axis.text.x = element_text(angle = 45)) +
387   theme_bw() +
388   xlab("Health expenditure per capita (% GDP) quantile") +
389   ylab("Life Expectancy (Years)") +
```

```
387 theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
388
389 ## Plot 3: Healthcare per capita (plus visualization of size per continent)
390 p3 <- ggplot(data = LE_data, aes(x = gdp_pc,
391                 y = life_expectancy))+
392     geom_point(aes(size = pop_total, color = Continent),
393                 show.legend = c(TRUE)) +
394     ggsci::scale_color_jama()+
395     scale_x_log10() +
396     theme_bw()+
397     ylab("Life Expectancy (Years)") +
398     xlab("GDP Per Capita ($)") +
399     guides(size = "none") +
400     theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
401
402 # Plot in a grid with cowplot::plot_grid() and also give titles
403 top <- cowplot::plot_grid(p1, p2, labels = c("A. Histogram of Life
404     Expectancy", "B. Life Expectancy by Healthcare Spending"),
405     hjust = -0.1)
406 bottom <- cowplot::plot_grid(p3, labels = c("C. Life Expectancy by
407     Population and GDP per capita"),
408     hjust = -0.1)
409 cowplot::plot_grid(top, bottom, ncol = 1)
410
411 #####
412 ##### Figure S1: Life expectancy by health expenditure per capita
413 ######
414 # Also plot life exp, health expenditure per capita for Figure S1 (to show
415 # we need to log-transform)
416 p1 <- ggplot(data = LE_data, aes(x = health_expend_pc,
417                 y = life_expectancy))+
418     geom_point() +
419     ggsci::scale_color_jama()+
420     theme_bw()+
421     ylab("Life Expectancy (Years)") +
422     xlab("Health spending per capita ($)") +
423     guides(size = "none") +
424     theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
425 p2 <- ggplot(data = LE_data, aes(x = health_expend_pc,
426                 y = life_expectancy))+
427     geom_point() +
428     ggsci::scale_color_jama()+
429     scale_x_log10() +
430     theme_bw()
```

```
428 ylab("Life Expectancy (Years)") +
429 xlab("log(Health spending per capita)") +
430 guides(size = "none") +
431 theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
432
433 cowplot::plot_grid(p1,p2,nrow = 1)
434
435 # Log-transform GDP per capita and health_expend_pc in dataset
436 LE_data$gdp_pc <- log(LE_data$gdp_pc)
437 LE_data$health_expend_pc <- log(LE_data$health_expend_pc)
438
439 #####
440 ## Section 1B: Missing Values -----
441 #####
442
443 # Subset variables with missingness < 50% (we identified these vars earlier)
444 LE_data1 <- LE_data[,names(LE_data) %in% low_missing_vars]
445 Country.Name <- LE_data$Country.Name
446
447 # Remove ID column to get final dataset for imputation
448 LE_data2 <- subset(LE_data1, select = -c(Country.Name))
449
450 # Now impute using using the mice package, which performs multiple
#      imputation by chained equations.
451 #?mice() to access documentation
452
453 # First, we create predictor matrix without LE, so LE isn't used to predict
#      covariates.
454 pred <- mice::make.predictorMatrix(LE_data2)
455 pred[, "life_expectancy"] <- 0 # never use it as a predictor
456
457 # Imputation using method predictive mean matching (pmm) and m = 10
#      iterations
458 imputations <- mice(LE_data2, seed = 98, predictorMatrix = pred, method =
#      "pmm", m = 10)
459
460 #####
461 ##### Figure 2: Plot of Imputations
462 #####
463
464 # stipplot() as we used in lab did not seem to work for some reason.
465 # So we plot the imputed and observed values for a subset of features
#      manually.
```

```
466 # For this, we just chose a somewhat random subset of features that covered
467 # different
468 #
469 # Countries missing life expectancy data
470 countries_missing_gdp_pc <- dplyr::filter(LE_data, is.na(gdp_pc))$Country.
471   Name
472 countries_missing_health_expend_pc <- dplyr::filter(LE_data, is.na(health_
473   expend_pc))$Country.Name
474 countries_missing_safe_drinking_water <- dplyr::filter(LE_data,
475                                         is.na(safe_drinking_
476   water))$Country.Name
477 countries_missing_infant_mortality_rate <- dplyr::filter(LE_data,
478                                         is.na(infant_
479   mortality_rate))$Country.Name
480
481 # We will only plot first 5 imputations so it is easy to see. Extract first
482 # 5 imputed datasets
483 vars <- c("life_expectancy", "adults_HIV", "health_expend_pc", "infant_
484   mortality_rate")
485 imp1 <- cbind(LE_data$Country.Name, "imp" = 1, complete(imputations, 1))
486 imp2 <- cbind(LE_data$Country.Name, "imp" = 2, complete(imputations, 2))
487 imp3 <- cbind(LE_data$Country.Name, "imp" = 3, complete(imputations, 3))
488 imp4 <- cbind(LE_data$Country.Name, "imp" = 4, complete(imputations, 4))
489 imp5 <- cbind(LE_data$Country.Name, "imp" = 5, complete(imputations, 5))
490 imps <- rbind(imp1, imp2, imp3, imp4, imp5)
491
492 # Re-add country name, which we had removed before imputing.
493 # Also define new variable specifying which points were imputed
494 names(imps)[1] <- "Country.Name"
495 imps <- imps %>%
496   dplyr::mutate(Type_gdp = ifelse(Country.Name %in% countries_missing_gdp_pc
497     ,
498       "Imputed", "Not Imputed"),
499       Type_safe_drinking = ifelse(Country.Name %in% countries_
500         missing_safe_drinking_water,
501           "Imputed", "Not Imputed"),
502       Type_health = ifelse(Country.Name %in% countries_missing_
503         health_expend_pc,
504           "Imputed", "Not Imputed"),
505       Type_infant = ifelse(Country.Name %in% countries_missing_
506         infant_mortality_rate,
507           "Imputed", "Not Imputed"))
```

```
499 p1 <- ggplot(arrange(imps, desc(Type_gdp)), aes(x = imp,
500                                         y = gdp_pc,
501                                         color = Type_gdp)) +
502   geom_point(size = 0.8, shape = 1,
503               show.legend = FALSE) +
504   theme_bw() +
505   scale_color_manual(values = c("#8B0000", "#FFFF00")) +
506   xlab("Imputation") +
507   ylab("log(GDP per capita)")
508
509 p3 <- ggplot(arrange(imps, desc(Type_safe_drinking)), aes(x = imp,
510                                         y = safe_drinking_water,
511                                         color = Type_safe_drinking)) +
512   geom_point(shape = 1, size = 0.8, show.legend = FALSE) +
513   theme_bw() +
514   scale_color_manual(values = c("#8B0000", "#FFFF00")) +
515   xlab("Imputation") +
516   ylab("Safe drinking water (%)")
517
518 p2 <- ggplot(arrange(imps, desc(Type_health)), aes(x = imp,
519                                         y = health_expend_pc,
520                                         color = Type_health)) +
521   geom_point(shape = 1, size = 0.8, show.legend = FALSE) +
522   theme_bw() +
523   scale_color_manual(values = c("#8B0000", "#FFFF00")) +
524   xlab("Imputation") +
525   ylab("log(Health expenditure per capita)")
526
527 p4 <- ggplot(arrange(imps, desc(Type_infant)), aes(x = imp,
528                                         y = infant_mortality_rate,
529                                         color = Type_infant)) +
530   geom_point(shape = 1, size = 0.8, show.legend = FALSE) +
531   theme(axis.title.y = element_text(size = 9)) +
532   theme_bw() +
533   scale_color_manual(values = c("#8B0000", "#FFFF00")) +
534   xlab("Imputation") +
535   ylab("Infant mortality rate (per 1,000 live births)")
536
537 cowplot::plot_grid(p1, p2, p3, p4, nrow = 2, hjust = -0.1)
538
539 #####
```

```
540 ##### Question 2: Dealing with Collinearity -----
541 #####
542
543 # Studying collinearity presence, which we can do on multiply imputed data
#       using
544 # the 'miceadds' package
545 #?miceadds
546
547 # First, we are calculating correlations by removing Continent (factor) and
#       life expectancy (response)
548 mi_corr <- micombine.cor(imputations, variables = c(3:21))
549
550 # This gives pairwise correlations in columns. We're now turning it into a
#       matrix (for plotting, analysis)
551 cor_df <- mi_corr %>%
552   dplyr::select(variable1, variable2, r) %>%
553   tidyr::pivot_wider(names_from = variable2, values_from = r)
554
555 # Rearranging columns and making diagonals 1 instead of NA
556 cor_df <- cor_df %>%
557   dplyr::relocate(electricity_access, .before = income_growth) %>%
558   dplyr::mutate_all(~replace(., is.na(.), 1))
559
560 # Lastly, move var names from col 1 into row names
561 cor_df <- as.data.frame(cor_df)
562 rownames(cor_df) <- cor_df$variable1
563 cor_df <- cor_df[,-1]
564
565 #####
566 ##### Figure 3: Correlation Plot
567 #####
568
569 # Set names for corplot printing
570 colnames(cor_df) <- c("Electricity access", "Income growth",
571                       "Income growth percent", "Children out of school",
572                       "Infant mortality rate", "Primary completion rate",
573                       "Real interest rate", "Percent pop. growth",
574                       "Pop. density", "Total pop.",
575                       "Health spending per capita",
576                       "Health spending perc. GDP",
577                       "Unemployment rate", "GDP growth",
578                       "GDP per capita", "Birth rate",
579                       "Adults HIV perc.", "Safe drinking water",
580                       "Education duration")
```

```
581 rownames(cor_df) <- colnames(cor_df)
582
583 # and plot
584 corplot(as.matrix(cor_df), tl.cex = 0.8, tl.col="black")
585
586
587 # Variance inflation factors (VIF). Find mean over imputations initially.
588 # Since we didn't learn a clear way to calculate VIF on multiply imputed
589 # datasets (and don't
590 # see one in 'miceadds'),
591 # we instead calculate it on each one individually and then report the mean
592 # VIF for each variable
593 # across the 10 datasets
594 vif_output <- list()
595 for (i in 1:10)
596 {
597   vif_output[[i]] <- faraway::vif(complete(imputations, i)[,-c(1,2)])
598 }
599
600 vif_output_means <- c()
601 for (i in 1:length(vif_output[[1]]))
602 {
603   val_vec <- c()
604   for (j in 1:10)
605   {
606     val_vec[j] <- vif_output[[j]][i]
607   }
608   vif_output_means[i] <- mean(val_vec)
609   names(vif_output_means)[i] <- names(vif_output[[1]][i])
610 }
611
612
613 # These are the variables with large VIFs and high pairwise correlations w/
614 # other vars
615 drop_vars <- c("Continent", "life_expectancy", "income_growth", "birth_rate", "electricity_access", "health_expend_pc")
616
617 # Now recalculate VIFs after dropping the drop_vars
618 vif_output_after <- list()
619 for (i in 1:10)
620 {
```

```
620 vif_output_after[[i]] <- faraway::vif(complete(imputations, i)[, !(names(
621 complete(imputations, i)) %in% drop_vars)])
622 }
623
624 vif_output_means_after <- c()
625 for (i in 1:length(vif_output_after[[1]]))
626 {
627   val_vec_after <- c()
628   for (j in 1:10)
629   {
630     val_vec_after[j] <- vif_output_after[[j]][i]
631   }
632   vif_output_means_after[i] <- mean(val_vec_after)
633   names(vif_output_means_after)[i] <- names(vif_output_after[[1]][i])
634 }
635 # Now we see lower means
636 vif_output_means_after
637
638 #####
639 ## Table S1
640 #####
641
642 # Somehow this table isn't printing the VIF after, but we manually added the
643 # values to the table
644 # from the above "vif_output_means_after"!
645 temp <- data.frame("Feature" = names(vif_output_means_after),
646                      "VIF After Dropping Vars" = vif_output_means_after)
647
648 tableS1_df <- data.frame("Feature" = colnames(cor_df),
649                           "VIF" = sprintf("%.2f", vif_output_means))
650 tableS1_df <- dplyr::left_join(tableS1_df, temp, by = "Feature")%>%
651   dplyr::mutate(VIF.After.Dropping.Vars =
652                 ifelse(is.na(VIF.After.Dropping.Vars),
653                       " ",
654                       sprintf("%.2f", VIF.After.Dropping.Vars)))
655
656 names(tableS1_df)[3] <- "VIF After Dropping Vars"
657
658
659 tableS1_table <- flextable(tableS1_df)
660
661 tableS1_table <- align(tableS1_table, j = 1, align = "left", part = "all")
```

```
662 tableS1_table <- align(tableS1_table, j = 2:3, align = "right", part = "all")
663 tableS1_table <- width(tableS1_table, width = c(1.4, 1.0, 1.0))
664 #table1_table <- height_all(table1_table, height = 0.30)
665 tableS1_table <- hrule(tableS1_table, rule = "at least", part = "all")
666 tableS1_table # Note: we had to edit this in word after to add correct "
667     "after" values
668 #####
669 ### Question 3: Making the Model -----
670 #####
671
672 #Defining the analysis we're performing. Only include subset of variables (
673     without drop_vars)
674 feature.selection1 <- expression(null.model1 <- lm(life_expectancy ~ 1),
675     model2 <- step(null.model1, scope = ~ income_growth_pc + children_out_of_
676         _school + infant_mortality_rate + primary_completion_rate + real_interest_
677         _rate + pop_growth_perc_annual + pop_dens + pop_total + health_expend_
678         perc_gdp + unemployment_rate + gdp_growth + gdp_pc + adults_HIV + safe_
679         drinking_water + ed_duration, direction = "forward"))
680
681 # Redefine list of imputations to remove life expectancy for countries
682 # with NA life expectancy originally. We don't want to use them here - only
683     in predictions later
684 countries_with_LE <- dplyr::filter(LE_data,
685                                         !(is.na(life_expectancy)))$Country.Name
686
687 imp_filt <- mice::filter(imputations,
688                           Country.Name %in% countries_with_LE)
689
690 # Using the with function to evaluate the above expression
691 step.fit <- with(imputations, feature.selection1)
692
693 # Extracting the number of times each variable appears in the m imputed
694     datasets
695 step.fit <- with(imp_filt, feature.selection1)
696 step.fit.models <- lapply(step.fit$analyses, formula)
697 step.fit.features <- lapply(step.fit.models, terms)
698 feature.frequency <- unlist(lapply(step.fit.features, labels))
699
700 #Stepwise feature selection with our imputed datasets yields the following
701     table
702 #showing the frequency of feature selection:
703 sort(table(feature.frequency), decreasing=TRUE)
```

```
696
697 # Pool results according to above model (using mice, like we did in lab)
698 fit <- with(imp_filt, lm(life_expectancy ~
699                         health_expend_perc_gdp + infant_mortality_rate +
700                         safe_drinking_water + pop_dens +
701                         gdp_pc))
702
703 fit$analyses %>% str(max.level = 1)
704
705 pooled_results <- mice::pool(fit)
706 summary(pooled_results)
707 pool.r.squared(pooled_results, adjusted = TRUE)
708
709 ######
710 ### Table 3: The final model
711 #####
712
713 m <- 10
714 df <- as.data.frame(summary(pooled_results))
715 table3_df <- c()
716 table3_df$Feature <- c("Intercept",
717                         "Health Expenditure (% GDP)",
718                         "Infant mortality rate",
719                         "% drinking safe water",
720                         "Population density (per sq. km)",
721                         "log(GDP per capita)")
722 table3_df$Estimate <- sprintf("%.4f", df$estimate)
723 table3_df$SE <- sprintf("%.5f", df$std.error)
724 table3_df$df <- sprintf("%.2f", df$df)
725 table3_df$P_val <- sprintf("%.4f", df$p.value)
726 table3_df$P_val <- ifelse(table3_df$P_val == "0.0000",
727                             "< 0.0001", table3_df$P_val)
728 table3_df$BestSubsetFrequency <- c(" ",
729                                     table(feature.frequency)[names(table(
730                                         feature.frequency)) == "health_expend_perc_gdp"] / m,
731                                     table(feature.frequency)[names(table(
732                                         feature.frequency)) == "infant_mortality_rate"] / m,
733                                     table(feature.frequency)[names(table(
734                                         feature.frequency)) == "safe_drinking_water"] / m,
735                                     table(feature.frequency)[names(table(
736                                         feature.frequency)) == "pop_dens"] / m,
737                                     table(feature.frequency)[names(table(
738                                         feature.frequency)) == "gdp_pc"] / m
739 )
```

```
734  
735 # recalculate VIF for these five vars only  
736 vif_output <- list()  
737 for (i in 1:10){  
  vif_output[[i]] <- faraway::vif(dplyr::select(  
    complete(imputations, i), health_expend_perc_gdp,  
    infant_mortality_rate,  
    safe_drinking_water, pop_dens, gdp_pc))  
}  
742  
743 vif_output_means2 <- c()  
744 for (i in 1:length(vif_output[[1]])){  
  val_vec <- c()  
  for (j in 1:10){  
    val_vec[j] <- vif_output[[j]][i]  
  }  
  vif_output_means2[i] <- mean(val_vec)  
  names(vif_output_means2)[i] <- names(vif_output[[1]][i])  
}  
751  
752 vif_output_means2  
753 table3_df$VIF <- c(" ",vif_output_means2)  
754 table3_df <- as.data.frame(table3_df)  
755  
756  
757 names(table3_df) <- c("Feature",  
  "Estimate",  
  "Standard Error",  
  "df",  
  "P-value",  
  "Proportion of Best Subsets Appeared in",  
  "VIF")  
764  
765 table3_table <- flextable(table3_df, col_keys = c("Feature",  
  "Estimate",  
  "Standard Error",  
  "df",  
  "P-value",  
  "Proportion of Best  
  Subsets Appeared in",  
  "VIF"))  
772  
773 table3_table <- align(table3_table,j =1, align = "left", part = "all")  
774 table3_table <- align(table3_table,j =2:7, align = "right", part = "all")  
775 table3_table <- width(table3_table, width = c(1.4,1.0, 1.1, 1.3, 1.0, 1.0,  
  1.0 ))
```

```
776 table3_table <- hruler(table3_table, rule = "at least", part = "all")
777 table3_table <- vline(table3_table, i = NULL, j = 5, part = "body",
778                         border = fp_border(color="black"))
779 table3_table
780
781 #####
782 ##### Figure 4: Residuals and model assumption check
783 #####
784
785
786 # Extract residuals to check model assumptions.
787 # We will do this for each imputed dataset individually and report averages.
788 # Should be ok with m = 10. A little higher and we might start to worry
    about CLT.
789 predict_df <- cbind(
790   predict(fit$analyses[[1]], newdata = complete(imp_filt, 1)),
791   predict(fit$analyses[[2]], newdata = complete(imp_filt, 2)),
792   predict(fit$analyses[[3]], newdata = complete(imp_filt, 3)),
793   predict(fit$analyses[[4]], newdata = complete(imp_filt, 4)),
794   predict(fit$analyses[[5]], newdata = complete(imp_filt, 5)),
795   predict(fit$analyses[[6]], newdata = complete(imp_filt, 6)),
796   predict(fit$analyses[[7]], newdata = complete(imp_filt, 7)),
797   predict(fit$analyses[[8]], newdata = complete(imp_filt, 8)),
798   predict(fit$analyses[[9]], newdata = complete(imp_filt, 9)),
799   predict(fit$analyses[[10]], newdata = complete(imp_filt, 10)))
800
801 observed_df <- cbind(
802   dplyr::select(complete(imp_filt, 1), life_expectancy),
803   dplyr::select(complete(imp_filt, 2), life_expectancy),
804   dplyr::select(complete(imp_filt, 3), life_expectancy),
805   dplyr::select(complete(imp_filt, 4), life_expectancy),
806   dplyr::select(complete(imp_filt, 5), life_expectancy),
807   dplyr::select(complete(imp_filt, 6), life_expectancy),
808   dplyr::select(complete(imp_filt, 7), life_expectancy),
809   dplyr::select(complete(imp_filt, 8), life_expectancy),
810   dplyr::select(complete(imp_filt, 9), life_expectancy),
811   dplyr::select(complete(imp_filt, 10), life_expectancy))
812
813 residual_df <- predict_df - observed_df # Note that residual_df has 10
    columns (one for each imp) and 198 rows (one row per obs).
814 names(residual_df) <- c(1:10)
815
816
```

```
817 # Now we find the mean for each obs and use that for residual diagnostic
  plots.
818 residual_means <- residual_df %>%
  dplyr::mutate(mean = rowMeans(residual_df))
819 residual_means <- residual_means[, "mean"]
820
821
822 # And do same for fit (for residual vs fit)
823 predicted_means <- as.data.frame(predict_df) %>%
  dplyr::mutate(mean = rowMeans(predict_df))
824 predicted_means <- predicted_means[, "mean"]
825
826
827 mean_df <- as.data.frame(cbind(predicted_means, residual_means))
828
829
830 # Plot
831
832 p1 <- ggplot(mean_df, aes(x=residual_means)) +
  geom_histogram() + ggsci::scale_colour_jama() +
  theme_bw() + theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm")) +
  xlab("Residuals") +
  ylab("Count")
833
834 p2 <- ggplot(mean_df, aes(sample = residual_means)) +
  stat_qq() + stat_qq_line() + theme_bw() +
  theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm")) +
  xlab("Theoretical") +
  ylab("Sample")
835
836 p3 <- ggplot(mean_df, aes(x=predicted_means, y = residual_means)) +
  geom_point() + theme_bw() +
  theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm")) +
  xlab("Fitted Value") +
  ylab("Residual")
837
838
839 # Plot's titles and format
840 top <- plot_grid(p1, p2, labels = c("A. Histogram of Residuals", "B. QQplot
  of Residuals"),
  hjust = -0.1)
841 bottom <- plot_grid(p3, labels = "C. Residuals vs. Fit", hjust = -0.1)
842 plot_grid(top, bottom, nrow = 2)
843
844
845 #####
846 #### Table S2: Prediction for other 19 countries
```

```
859 #####  
860  
861 # Generate predictions for the 19 countries initially missing LE and print  
# results in table  
862  
863 # First re-add Country Name variable to each imputed dataset  
864 dataList <- list(cbind(LE_data$Country.Name, complete(imputations,1)),  
865                 cbind(LE_data$Country.Name, complete(imputations,2)),  
866                 cbind(LE_data$Country.Name, complete(imputations,3)),  
867                 cbind(LE_data$Country.Name, complete(imputations,4)),  
868                 cbind(LE_data$Country.Name, complete(imputations,5)),  
869                 cbind(LE_data$Country.Name, complete(imputations,6)),  
870                 cbind(LE_data$Country.Name, complete(imputations,7)),  
871                 cbind(LE_data$Country.Name, complete(imputations,8)),  
872                 cbind(LE_data$Country.Name, complete(imputations,9)),  
873                 cbind(LE_data$Country.Name, complete(imputations,10)))  
874  
875 # And filter to only keep  
876 dataList2 <- list()  
877 for (i in 1:length(dataList)){  
878   names(dataList[[i]])[1] <- "Country.Name"  
879   dataList2[[i]] <- dplyr::filter(dataList[[i]],  
880                                     Country.Name %in% countries_missing_LE)  
881 }  
882  
883 # Now predict over each one using coefficients from pooled model  
884 # We aren't aware of an easy way to use 'predict()' on a model output by  
# mice,  
885 # so we perform predictions for estimates manually using the pooled  
# coefficient estimates.  
886 coefs <- pooled_results$pooled$estimate  
887 pred_list <- list()  
888  
889 for (i in 1:length(dataList2)){  
890   df <- dataList2[[i]]  
891   pred_list[[i]] <- coefs[1] + coefs[2]*df$health_expend_perc_gdp +  
892   coefs[3]*df$infant_mortality_rate + coefs[4]*df$safe_drinking_water+  
893   coefs[5]*df$pop_dens + coefs[6]*df$gdp_pc  
894 }  
895  
896 # Take mean of prediction for each country across imputed datasets  
897 final_preds <- cbind(dataList2[[1]]$Country.Name, bind_cols(pred_list))  
898  
899 names(final_preds) = c("Country.Name",c("A1","A2",
```

```
900                     "A3", "A4",
901                     "A5", "A6",
902                     "A7", "A8",
903                     "A9", "A10"))
904
905 final_preds <- dplyr::mutate(final_preds,
906                               mean_pred =
907                                 rowMeans(dplyr::select(final_preds, starts_
908                               with("A"))))
909
910 # Now plot
911 tableS2_df <- dplyr::select(final_preds, Country.Name, mean_pred) %>%
912   dplyr::mutate(mean_pred = sprintf("%.2f", mean_pred))
913
914 names(tableS2_df) <- c("Country", "Predicted life expectancy")
915
916 tableS2_table <- flextable(tableS2_df,
917                             col_keys = c("Country",
918                                         "Predicted life expectancy"))
919
920 tableS2_table <- align(tableS2_table, j = 1, align = "left", part = "all")
921 tableS2_table <- align(tableS2_table, j = 2, align = "right", part = "all")
922 tableS2_table <- width(tableS2_table, width = c(1.4, 1.0))
923 #table1_table <- height_all(table1_table, height = 0.30)
924 tableS2_table <- hrule(tableS2_table, rule = "at least", part = "all")
925 tableS2_table
926
927 #####
928 ### Question 4: Experimental Design -----
929 #####
930
931 # Performing ANOVA III
932 anova_LE <- aov(life_expectancy ~ as.factor(Continent), data=complete(
933   imputations))
934 summary(anova_LE)
935 Anova(anova_LE, type = "III")
936
937 # This specific section of code wasn't from lab.
938 # We used this R-bloggers post for reference [5] here:
939 # https://www.r-bloggers.com/2021/07/how-to-perform-ancova-in-r/
940
941 # We just pick complete(imputations, 1) as our single imputation
942 anova_data <- complete(imputations, 1)
```

```
942 anova_data$Continent <- as.factor(anova_data$Continent)
943
944 # Performing ANCOVA (according to slides from Lecture 10)
945 le.aov <- aov(life_expectancy ~ log(gdp_pc) + Continent, data = anova_data)
946 summary(le.aov)
947
948 # Post-hoc test
949 tukey.le <- glht(le.aov, linfct = mcp(Continent = "Tukey"))
950
951 # view a summary of the post-hoc comparisons
952 summary(tukey.le)
953
954 # Note: we also tried re-running above code using complete(imputations, i)
955 # With i = 2,3,...,10 as sensitivity analysis and got similar results
956
957 #Before proceeding with the statistical tests, we check the assumptions of
958 #ANCOVA which are (show diagnostic plots in appendix):
959
960 #####
961 ### Figure 5:
962 #####
963
964 # Plot life expectancy by % of GDP spent on healthcare
965 p1 <- ggplot(anova_data, aes(x = log(gdp_pc), y = life_expectancy, color =
966   Continent)) +
967   theme_bw() +
968   geom_point(alpha = 0.5, size = 0.8) + xlab("log(GDP per capita in
969   thousands)") +
970   ylab("Life expectancy") + ggsci::scale_color_jama() + geom_smooth(method =
971     "lm", se = FALSE)
972
973 # Figure showing life expectancy by continent (with singly imputed dataset
974 # that we will use)
975 p2 <- ggplot(complete(imputations, 1), aes(x = Continent, y = life_
976   expectancy, color = Continent)) +
977   geom_boxplot(show.legend = FALSE) + ggsci::scale_color_jama() + theme_bw()
978   + theme(axis.text.x = element_text(angle = -45),
979
980       plot.margin = unit(c(1,0.5,0.5,0.5), "cm")) + xlab("Continent")
981   + ylab("Life Expectancy (Years)")
982
983 plot_grid(p1, p2, labels = c("A", "B"))
984
985 # Test for homogeneity of variances
```

```
978 car::leveneTest(anova_data$life_expectancy~anova_data$Continent)
979
980 # Analysing residual's normality assumptions
981 anova_data$residuals <- le.aov$residuals
982
983
984 #####
985 ### Figure S2
986 #####
987
988 p1 <- ggplot(anova_data, aes(residuals)) +
989   geom_histogram(color = "#000000", fill = "#0099F8") +
990   theme_bw() +
991   xlab("Residual") +
992   ylab("Frequency") +
993   theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm"))
994
995 # Plot 2
996 p2 <- ggplot(anova_data, aes(sample = residuals)) +
997   stat_qq() +
998   stat_qq_line() +
999   theme_bw() +
1000   theme(plot.margin = unit(c(1,0.5,0.5,0.5), "cm")) +
1001   xlab("Theoretical") +
1002   ylab("Sample")
1003
1004 plot_grid(p1, p2, labels = c("A.", "B"))
1005
1006 # Shapiro test
1007 shapiro.test(anova_data$residuals)
```

8 Statement on contribution of group members

Students [2200634] and [2200909] participated equally in this project. We worked together in-person to create the bones of the code, write much of the report, and prepare the slides. For the parts that we did split up, [2200909] worked on the figures/tables and model fine-tuning while [2200634] worked on the continent model, documentation, and assembling the slides. We still reviewed and worked on the final version together for each of these tasks. One more student (unknown registration number) was also assigned to our group but did not contribute at all because they stopped responding to our messages or contacting us.