

Your project title

Last!: Parker Dingman, Ethan Donecoff, Karam Oubari, Pei Yi Zhuo

2021-04-17

Your written report goes here! Before you submit, make sure your code chunks are turned off with `echo = FALSE` and there are no warnings or messages with `warning = FALSE` and `message = FALSE`

Introduction and Data

By telling us the average age of death in a population, life expectancy is a key metric for understanding a country's health. According to Max Roser, Esteban Ortiz-Ospina and Hannah Ritchie from Our World in Data, "Broader than the narrow metric of the infant and child mortality, which focus solely at mortality at a young age, life expectancy captures the mortality along the entire life course." [1] Over the course of history, life expectancy has risen dramatically. It is estimated that in pre-modern times, life expectancy worldwide was only about 30 years. Since the Industrial Revolution in the 18th and 19th centuries, many countries had huge gains in this number. And since the beginning of the 20th century, global average life expectancy has risen to about 70 years. However, there remain huge inequalities in this number. Currently (as of 2019), the Central African Republic has the lowest life expectancy of 53 years while Japan has the highest with 83.

In addition to the more obvious health-related connections to life expectancy, numerous pieces of academic literature have delved into the non-medical factors behind life expectancy. A major example is a longitudinal study conducted by Charles Lin, Eugene Rogot, Norman Johnson, Paul Sorlie, and Elizabeth Arias, which examined life expectancy by socioeconomic factors. [2] Academic literature such as this provides us with motivations to examine this topic on an international level, looking at various health-related and non-health-related factors that connect to life expectancy.

In terms of initial hypotheses of model selection, we expect that the strongest predictors of life expectancy will be `Adult Mortality`, `infant deaths`, and `GDP`. We also predict that countries that have `status` equal to "Developed" will have higher life expectancy than those that have `status` equal to "Developing". We also predict that, on average, life expectancy will have increased internationally in the 15 years that are documented in the data set.

Our selected data set is comprised of information that others gleaned from the websites of the World Health Organization and the United Nations. Every entry describes the health, social, and economic conditions for one of 193 countries in a given year from 2000 to 2015. Additional information can be found on <https://www.kaggle.com/kumarajarshi/life-expectancy-who>.

We also combined this data set with another data set downloaded from <https://www.kaggle.com/fernandol/countries-of-the-world>.

Exploratory Data Analysis

First, we will look at some summary statistics of our response variable, 'Life expectancy'. Though it is possible to analyze the entire dataset over all years, this creates difficulties with creating models. As one might expect, the average global life expectancy increased from 2000 to 2015. This relationship might make it unclear whether a rise in life expectancy is explained by our predictor variables or if it is simply due to human development over time. As a result, we will only use data from one year to perform our analysis.

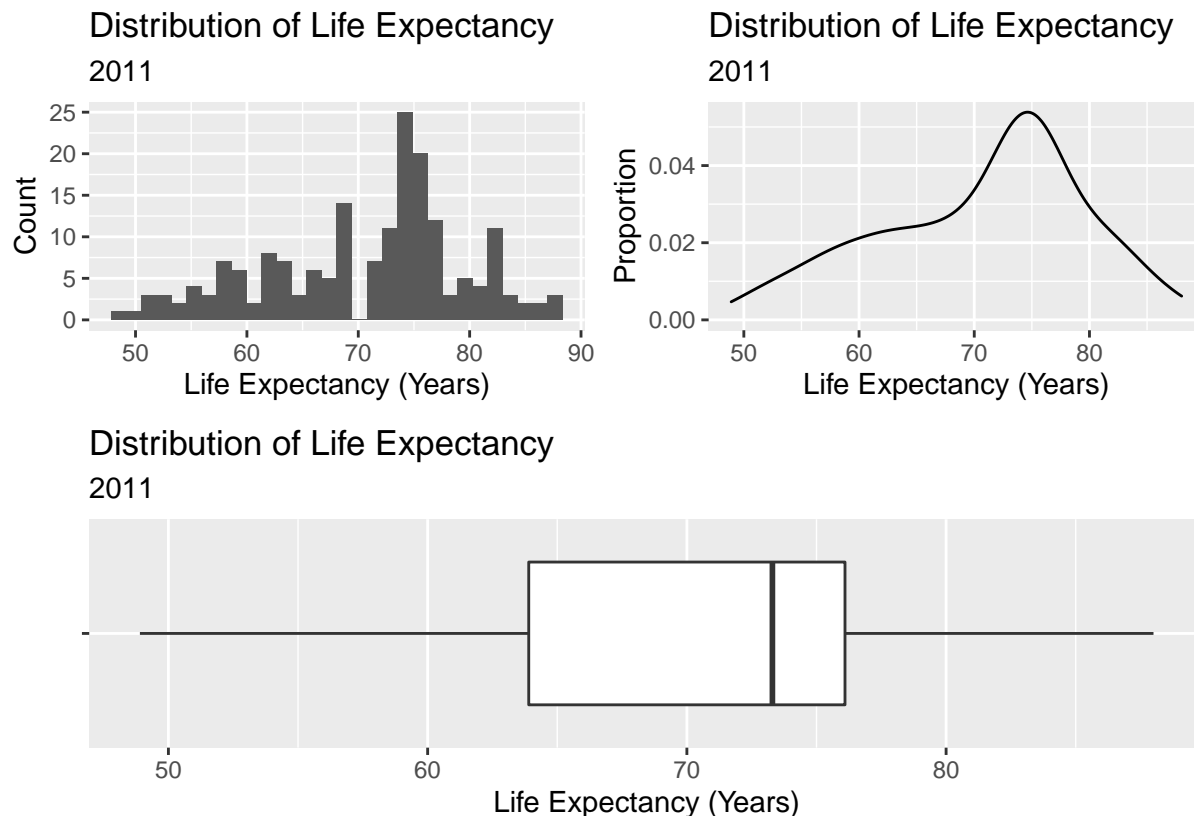
We will only analyze life expectancy for the year 2011, because it offered the most complete data (the most data without missing values).

Table 1: Summary Statistics of Response Variable, Life Expectancy

min	max	range	mean	median	Q1	Q3	iqr	sd
48.9	88	39.1	70.654	73.3	63.9	76.1	12.2	8.925

These summary statistics give a rough idea of the distribution of the response variable. The median life expectancy (~73.3 yrs) is almost 3 years larger than the mean (~70.7 yrs). Additionally, the median is closer to the third quartile than the first quartile. This suggests that life expectancy may be left-skewed, which we will evaluate further with visualizations.

Now, we will look at some summary visualizations of life expectancy.



Here we see that life expectancy is left skewed with a center just over 70 years. Life expectancy ranges from about 50 years to 90 years.

Methodology

Model Selection

In this project we plan on using various regression analysis methods including, but not limited to, multiple linear regression, statistical inference, analysis of variance, and model selection in an attempt to understand country and region-level life expectancy, as well as the health, social, and economic relationships behind this number.

The response variable we will be using in this project is **Life expectancy**. This is a measure of the average age of death in a year for the given country during that year. In terms of regression model technique, we will be using multiple linear regression because our response variable is quantitative.

The goal of our analysis is to calculate the most precise prediction of the response variable (life expectancy). To do this we will do backwards selection using AIC which means dropping variables one at a time that are deemed irrelevant based on AIC.

Table 2: Full Model With All Possible Predictor Variables

term	estimate	std.error	statistic	p.value
(Intercept)	49.126	3.094	15.876	0.000
RegionBALTICS	-1.206	2.052	-0.588	0.558
RegionC.W. OF IND. STATES	-1.763	1.473	-1.197	0.234
RegionEASTERN EUROPE	0.980	1.700	0.577	0.565
RegionLATIN AMER. & CARIB	1.931	1.336	1.445	0.151
RegionNEAR EAST	1.333	1.678	0.795	0.429
RegionNORTHERN AFRICA	1.399	1.855	0.754	0.452
RegionNORTHERN AMERICA	5.640	3.376	1.671	0.098
RegionOCEANIA	0.458	1.641	0.279	0.781
RegionSUB-SAHARAN AFRICA	-1.361	1.225	-1.111	0.269
RegionWESTERN EUROPE	6.388	1.773	3.603	0.000
StatusDeveloping	1.105	1.243	0.888	0.376
Adult Mortality	-0.014	0.003	-4.443	0.000
infant deaths	0.091	0.041	2.249	0.027
Alcohol	0.008	0.106	0.079	0.937
percentage expenditure	0.001	0.001	1.204	0.232
Hepatitis B	0.050	0.019	2.623	0.010
Measles	0.000	0.000	-1.125	0.263
BMI	0.007	0.021	0.324	0.747
under-five deaths	-0.064	0.029	-2.184	0.031
Polio	-0.020	0.014	-1.440	0.153
Total expenditure	0.141	0.117	1.206	0.231
Diphtheria	-0.036	0.024	-1.520	0.132
HIV/AIDS	-0.529	0.121	-4.355	0.000
GDP	0.000	0.000	-1.036	0.303
Population	0.000	0.000	-1.292	0.199
thinness 1-19 years	0.199	0.272	0.731	0.466
thinness 5-9 years	-0.099	0.262	-0.378	0.706
Income composition of resources	43.132	6.029	7.154	0.000
Schooling	-0.548	0.250	-2.192	0.031

Table 3: Backward Model Selection with AIC

term	estimate	std.error	statistic	p.value
(Intercept)	51.884	2.331	22.259	0.000
RegionBALTICS	-2.391	1.787	-1.338	0.184
RegionC.W. OF IND. STATES	-2.564	1.169	-2.192	0.030
RegionEASTERN EUROPE	-0.093	1.304	-0.071	0.943
RegionLATIN AMER. & CARIB	1.148	1.018	1.128	0.262
RegionNEAR EAST	0.814	1.414	0.575	0.566
RegionNORTHERN AFRICA	1.115	1.691	0.659	0.511
RegionNORTHERN AMERICA	4.549	3.107	1.464	0.146
RegionOCEANIA	0.057	1.201	0.047	0.962
RegionSUB-SAHARAN AFRICA	-1.888	1.057	-1.785	0.077
RegionWESTERN EUROPE	5.154	1.306	3.946	0.000
Adult Mortality	-0.014	0.003	-4.517	0.000
infant deaths	0.053	0.031	1.674	0.097
Hepatitis B	0.036	0.018	2.019	0.046
Measles	0.000	0.000	-1.442	0.152
under-five deaths	-0.037	0.023	-1.611	0.110
Total expenditure	0.171	0.106	1.612	0.110
Diphtheria	-0.040	0.022	-1.773	0.079
HIV/AIDS	-0.560	0.118	-4.760	0.000
Income composition of resources	42.054	5.296	7.941	0.000
Schooling	-0.564	0.239	-2.355	0.020

Table 4: Backward Model Selection with BIC

term	estimate	std.error	statistic	p.value
(Intercept)	50.004	1.752	28.533	0.000
Adult Mortality	-0.018	0.003	-5.823	0.000
percentage expenditure	0.000	0.000	2.864	0.005
HIV/AIDS	-0.568	0.123	-4.626	0.000
Income composition of resources	46.054	4.843	9.510	0.000
Schooling	-0.558	0.244	-2.292	0.024

```
## # A tibble: 1 x 1
##   adj.r.squared
##   <dbl>
## 1      0.918

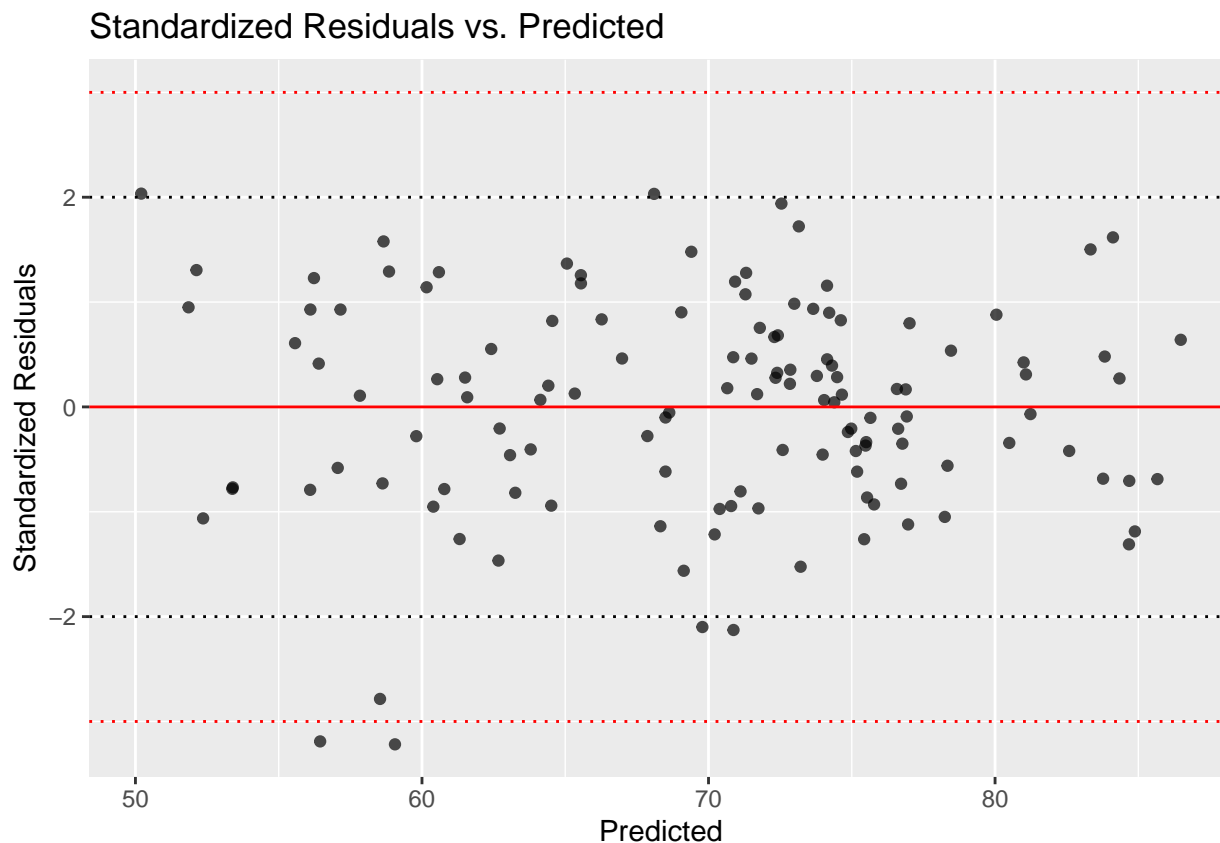
## # A tibble: 1 x 1
##   adj.r.squared
##   <dbl>
## 1      0.895
```

Based on adjusted R^2 , we prefer the model we found through backward selection using AIC. This model results in a higher adjusted R^2 which means that the model's predictor variables explain a higher proportion of the variation in the response (Life expectancy).

```
## # A tibble: 1 x 1
##   r.squared
##   <dbl>
## 1      0.930
```

93.03% of the variation in the Life expectancy is explained by the regression model containing Region, Adult Mortality, infant deaths, Hepatitis B, Measles, under-five deaths, Total expenditure, Diphtheria, HIV/AIDS, Income composition of resources and Schooling.

Model Diagnostics



There are a few moderate outliers and two severe outliers. Will look into these later.

```
## [1] 0.3230769
```

0.3230769 is the leverage threshold.

```
## # A tibble: 15 x 2
##   obs_num .hat
##   <int> <dbl>
## 1      3 0.345
## 2     25 1
## 3     40 0.335
## 4     43 0.328
## 5     54 0.756
## 6     55 0.410
## 7     65 0.339
## 8     69 0.346
## 9     82 0.348
## 10    84 0.339
## 11    87 0.403
## 12    90 0.711
## 13    91 0.373
## 14   121 0.340
```

```
## 15      125 0.375
## # A tibble: 0 x 2
## # ... with 2 variables: obs_num <int>, .cooks_d <dbl>
```

The points in the first table are points with high leverage ($\text{.hat} > \text{leverage_threshold} = 0.3230769$) while the points in the second table are influential points ($\text{.cooks_d} > 0.5$). 15 observations have high leverage and 0 observations are influential.

names	x
RegionBALTICS	1.408
RegionC.W. OF IND. STATES	1.899
RegionEASTERN EUROPE	1.920
RegionLATIN AMER. & CARIB	2.742
RegionNEAR EAST	1.721
RegionNORTHERN AFRICA	1.260
RegionNORTHERN AMERICA	1.441
RegionOCEANIA	1.630
RegionSUB-SAHARAN AFRICA	4.657
RegionWESTERN EUROPE	3.002
Adult Mortality	2.254
infant deaths	244.366
Hepatitis B	2.972
Measles	3.156
under-five deaths	244.039
Total expenditure	1.328
Diphtheria	2.777
HIV/AIDS	2.016
Income composition of resources	13.166
Schooling	8.937

Variables with a $VIF > 10$ will have issues with multicollinearity. **infant deaths** and **under-five deaths** are clearly highly correlated (this makes a lot of sense in the context of the data). **Income composition of resources** also seems to have issues with multicollinearity but we probably need to look into this more since it's not clear which other variable it is correlated with.

We should try models without either **infant deaths** and **under-five deaths** and then use model comparison techniques to decide on which of these two variables should be removed.

Results

Sources

- [1] <https://ourworldindata.org/life-expectancy>
- [2] <https://europepmc.org/article/med/12785422/reload=0#impact>