# STATISTICAL ANALYSIS TO UNDERSTAND LIFE EXPECTANCY

Abdallah Zaher, Cristina Mousinho, Gabriel Ravi, Nicolae-Radu

NOVA Information Management School and Universidade Nova de Lisboa

## Introduction

Life expectancy is a statistical measure of the average time an organism is expected to live, based on the year of its birth, its current age and other demographic factors including gender. Life expectancy has been growing for the last 10 decades, but not at the same rate in all countries.

## Objectives

This study is going to analyse 17 sustainable variables for 195 countries to understand which actions and cares can prolong the Life Expectancy in the world. The techniques that will be applied are **Stepwise Regression** and **Principal Components analysis**. This study contains:
- Descriptive analysis of all variables
- Analysis of the most important effects on Life Expectancy

## Data Exploration

The data was obtained from Open Access DataWorldBank: https://data.worldbank.org/ and hosted on: https://github.com/GRaviSantos79/StatsForDS.
The MicroWorldBankdata Library is a collection of data sets from the World Bank and other international, regional and national organizations. This data set contains information about the most important indicators of sustainability of the world.
This analysis consists of a longitudinal collection of 20 indicators of sustainability about 195 countries and response variables, Life Expectancy.
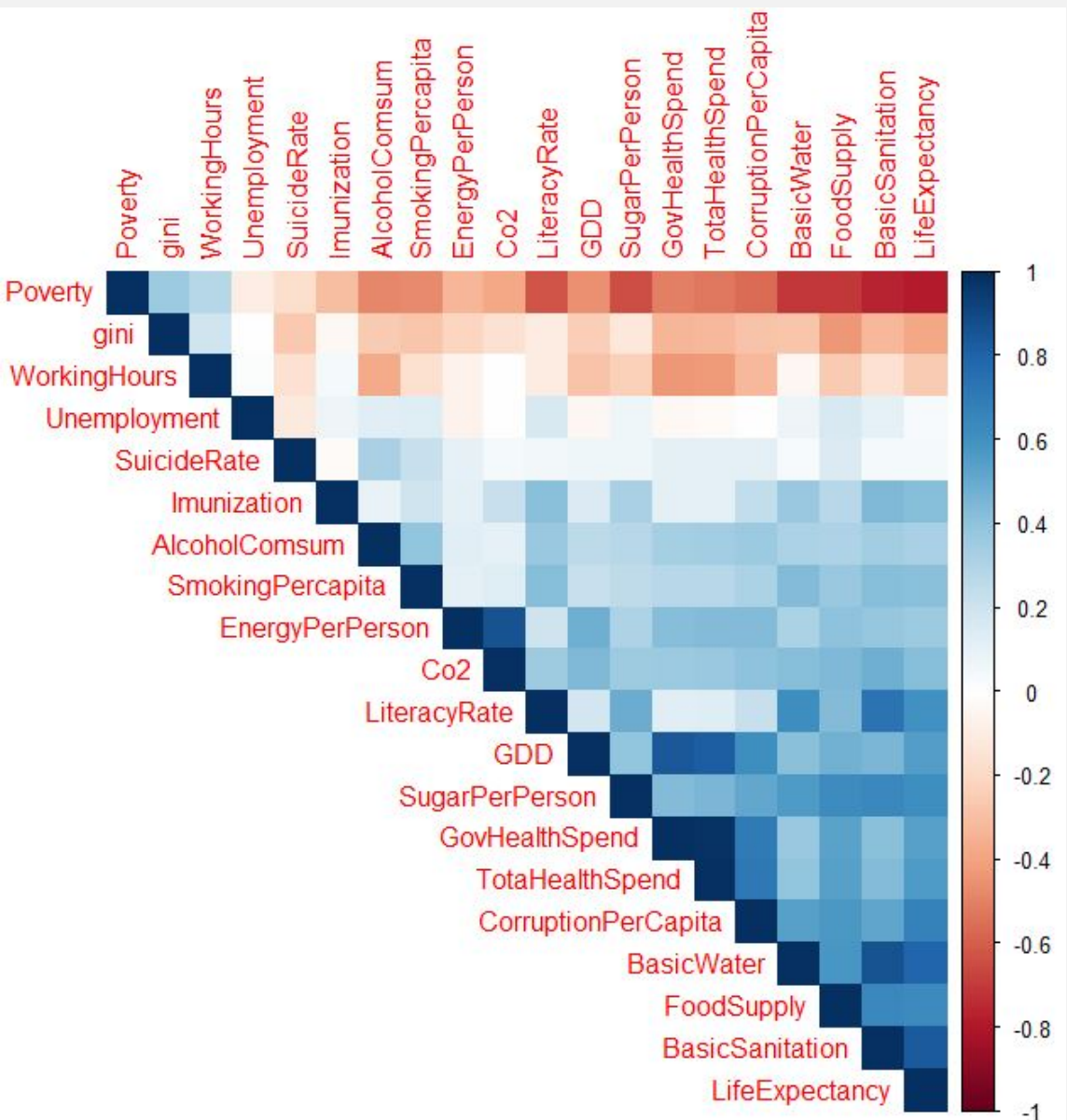The image below contains the descriptive analysis of all varaibles that may hav impact on Life Expectancy.



## Data Exploration

Check the summary of the variables bellow:

| BasicWater | % of Sugar/person |
|---|---|
| CO2 | CO2 emissions |
| GINI | Dist. of whealth |
| BasicSanitation | % Sanitation Available |
| Energy | % of Energy/person |
| Sugar | % of Water/person |
| GDD | Growing Degree Day |
| Smoking | % pop that smokes |
| Imunization | % pop immunized |
| Suicide | Number of Suicides/pop |
| WorkingHours | Working-Hours/person |
| Literacy | Literacy/person |
| Corruption | % of corruption |
| GovHealth | Gov Spend with health |
| FoodSupply | % of Food/person |
| Unemployment | % of Unemployment |
| TotalHealth | Pop Spend with health |
| Alcohol | % of Alcohol/person |

Here you can check and analyse a correlation graph between the variables.



Basic sanitation, food supply, basic water, poverty, corruption and GINI (negative correlation) are the variables that contain a higher correlation with Life Expectancy.
We will analyse these variables on further analysis because they are good candidates for the final model.
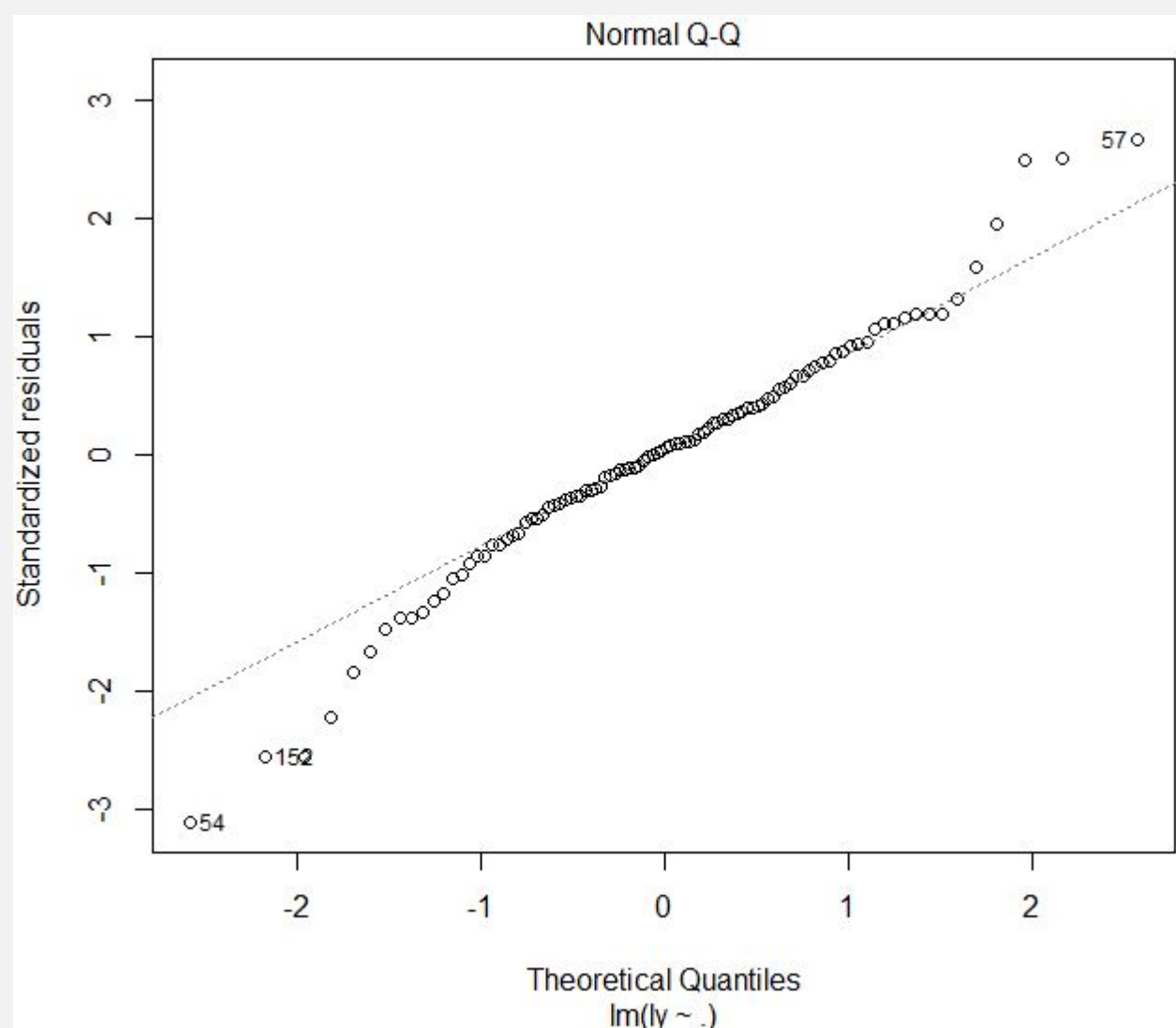
## Methodology

Life Expectancy will be converted to a logarithm scale because there are non-linear relationships in the data set and the chances of producing errors may also be skewed negatively. There is 5% of missing data. If we delete the countries that have NaN values, the data set will lose its quality. The missing values were filled by the median, as it is the technique that doesn't interfere with the results.

```
# Analysing and transforming the main variables
summary(data1$LifeExpectancy)
y = data1$LifeExpectancy
ly = log(y)
# Checking and treating the NA values for the training set
apply(Z, 2, function(X) sum(is.na(X)) )
X = ifelse(is.na(X), median(X, na.rm = TRUE), X)
```

The training set size is equal to 30% of the population. We made sure that the training set contains heterogeneous observations.

```
train = sample(1:nrow(data1), size = 58, replace = FALSE)
train;test = (-train) ;y.test = y[test] ;Xtreino = X[train,] ;xtest = x[test,]
```

The next step is to check the normality on the training set.



When we analyse the standardized residuals in each theoretical quantile, we can see that there are values outside the qqline, so we can't assume normality on the data set.
The heteroscedasticity was tested using the Breusch Pagan Test, at 95% of significance level, we can say that the null hyphotesis is not rejected, in other words, we don't have a problem with heteroscedasticity.
It is also important to check if any predictors are correlated with each other to avoid multicolineraity and guarantee that the model will not suffer of instability.

```
# MultiColinearity for the numerical variables
COR <- cor(data1)
Multt = findCorrelation(COR, cutoff = 0.75)
data1 = data1[,-Multt]
```

## Implementation

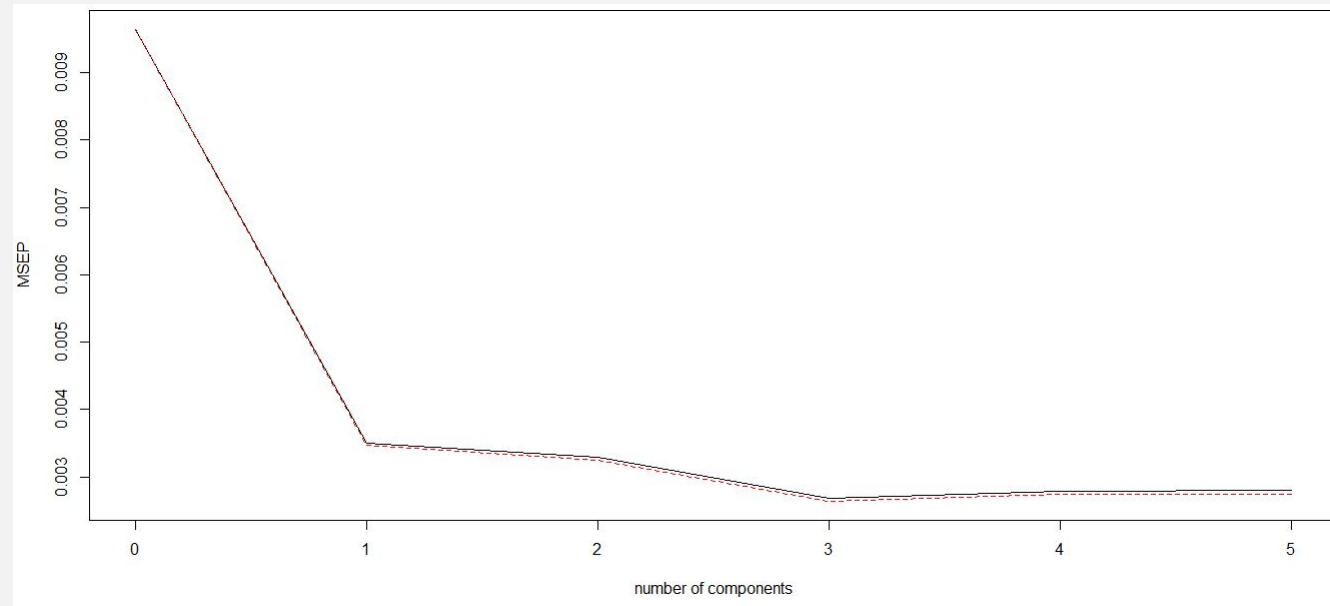1. The **Stepwise Regressions** (Backward, Forward and both) will be applied.

|  | $R^2$ | Nº Vars |
|---|---|---|
| **Backward** | 0.82 | 4 |
| **Foward** | 0.81 | 4 |
| **Both** | 0.83 | 5 |

When we consider the $R^2$, we see that the Both stepwise regression gives better estimates. After the analysis, considering 95% of confidence, the main variables are Poverty, Basic Sanitation, Basic Water, GINI and Alcohol Consumption. And the final model is:

$$LifeExpec = 4.2 - 0.03 Poverty + 0.04 Sanitation$$
$$+0.03 BasicWater - 0.02 GINI - 0.01 AlcoholCommsum \quad (1)$$

2. **Principal Components Analysis**, we will use the Mean squared prediction error to determine how many components.



To performe Principal Component Analysis, we will use the mean squaredprediction error, in order to determine how many components. In this case we chose 2 principal components, who explain 72% of the total variance.

## Conclusion

After the analysis, we can say that the most important variables to increase Life Expectancy are Poverty, BasicSanitation, Basic Water, GINI and alcoholConsum; in other words, these variables have the largest impact on life expectancy.
For example, in the final model(1), when a government increases basic water access by 3%, the Life Expectancy increases 1%.

## References

- Bendel, R.B. and Afifi, A.A., 1977. Comparison of stopping rules in