

# Life Expectancy Data: Analysis and Evaluation

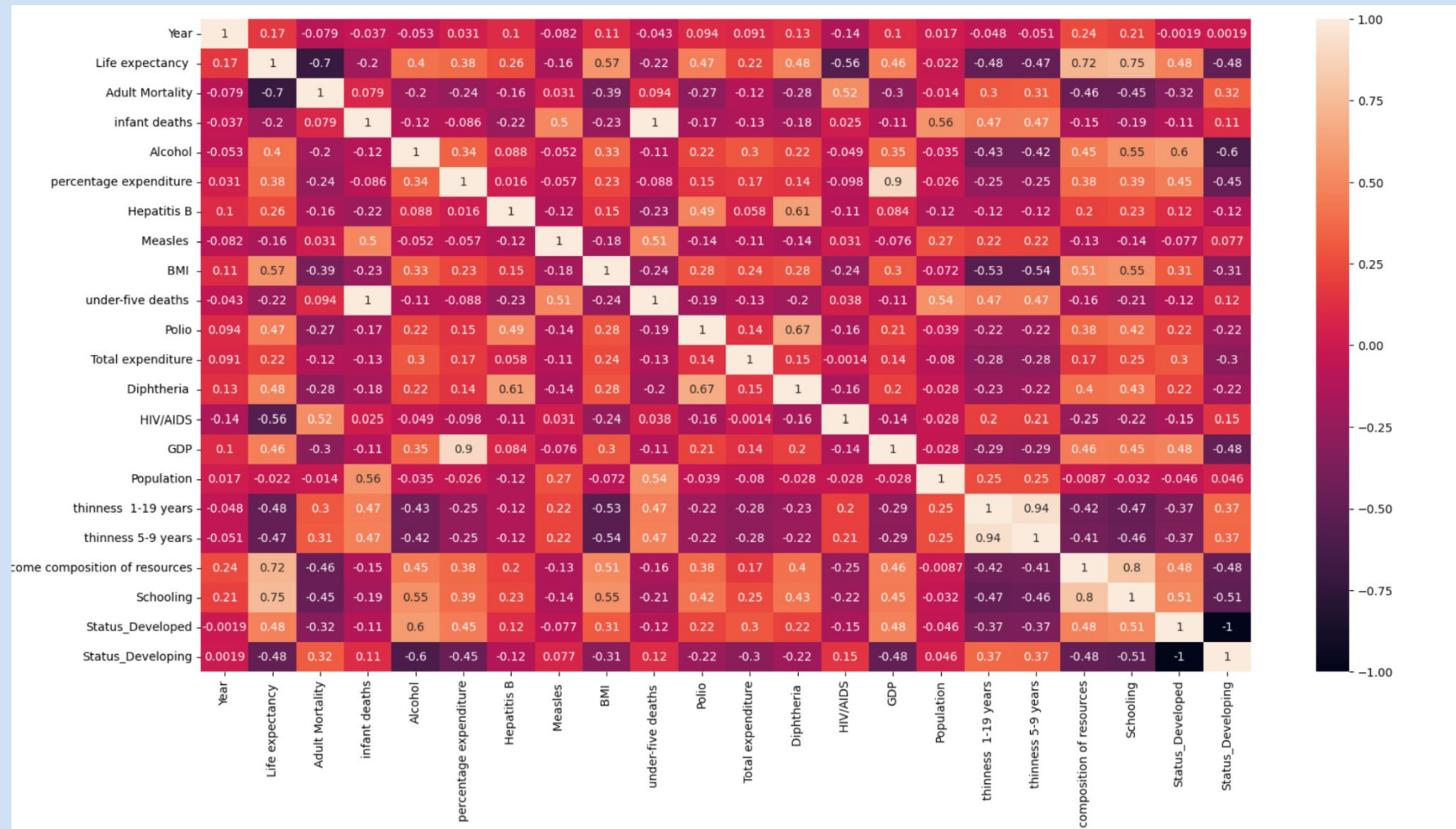
Harsh Panwar, Carlos Bonilla, Talha Choudhry, Nishtha Dandriyal  
CS301

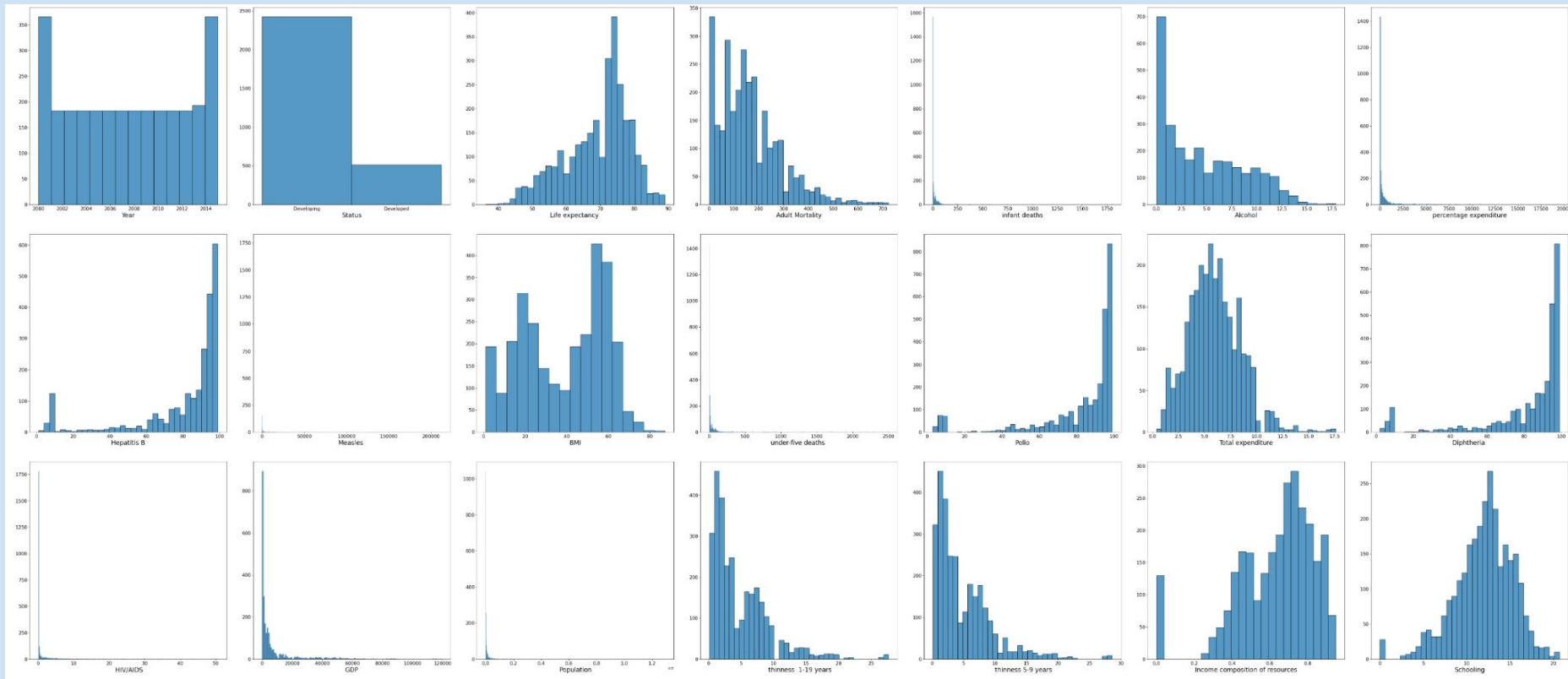
# Data Introduction + Pre-processing

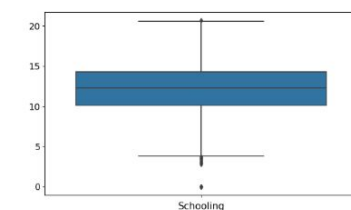
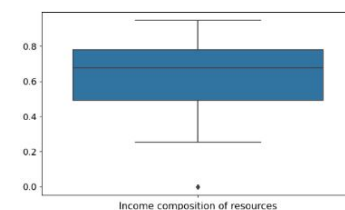
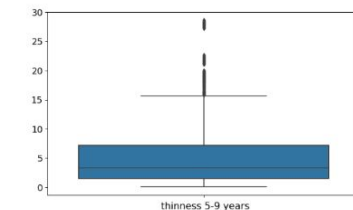
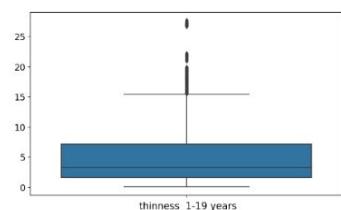
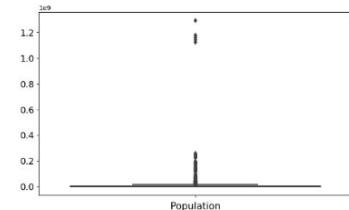
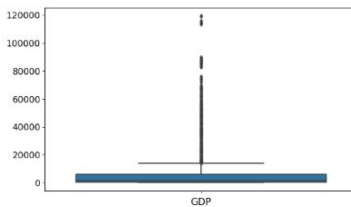
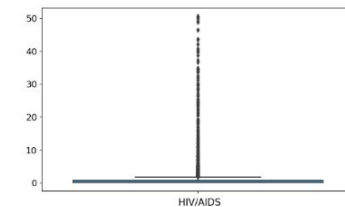
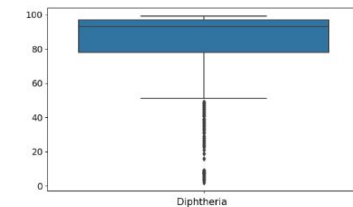
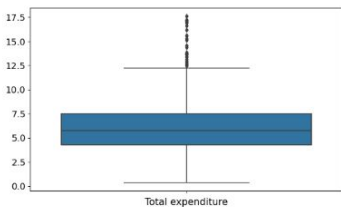
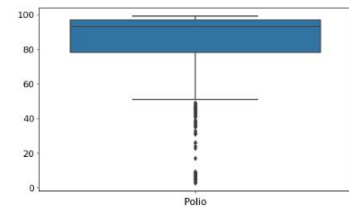
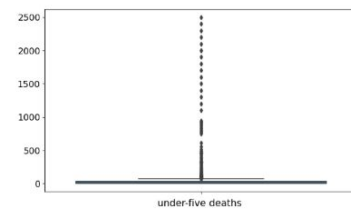
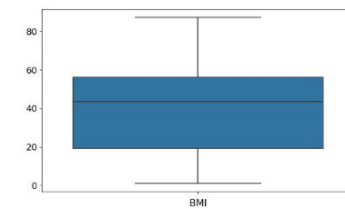
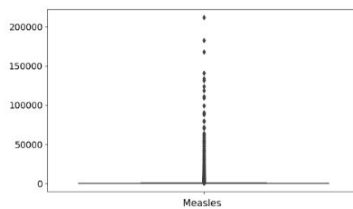
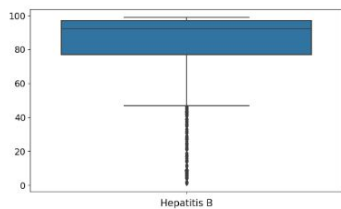
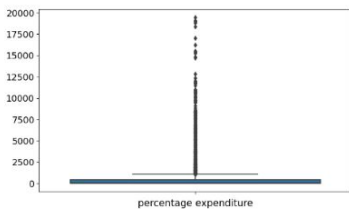
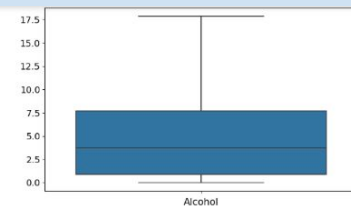
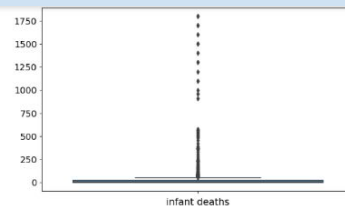
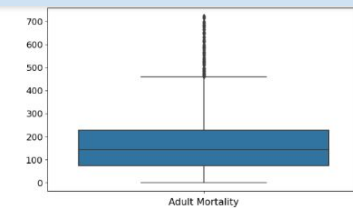
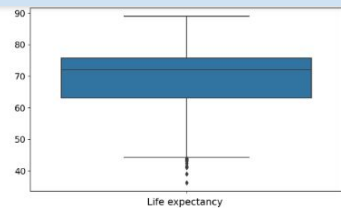
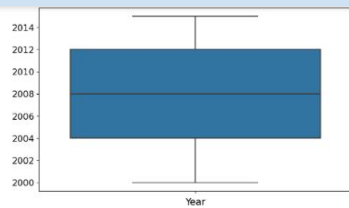
- Kaggle — sourced from World Health Organization
- 22 features
  - 20 Numerical
  - 2 Categorical
- One hot encoding for categorical features
- Dropped rows with any null values prior to training

# Dataset Analysis

Important Features	Significance
Adult Mortality	Probability of dying between 15 and 60 years per 1000 population
Alcohol	Alcohol, recorded per capita (15+) consumption (in litres of pure alcohol)
BMI	Average Body Mass Index of the entire population
HIV/AIDS	Deaths per 1,000 live births HIV/AIDS (0-4 years)
GDP	Gross Domestic Product per capita (in USD)
Income	Average amount of money earned
Schooling	Sum of the age specific enrollment rates for levels of education
Development status	Developed or developing?
Thinness	Prevalence of thinness among children and adolescents for Age 10 to 19 (%)







# Training Models

- Multiple Linear Regression
  - Optimization: Gradient Descent
  - Cost Function: Mean Square Error
- Decision Tree
  - Max Depth = 4
  - Optimization: CART
  - Cost Function: Mean Square Error

# Model Evaluation

- Multiple Linear Regression

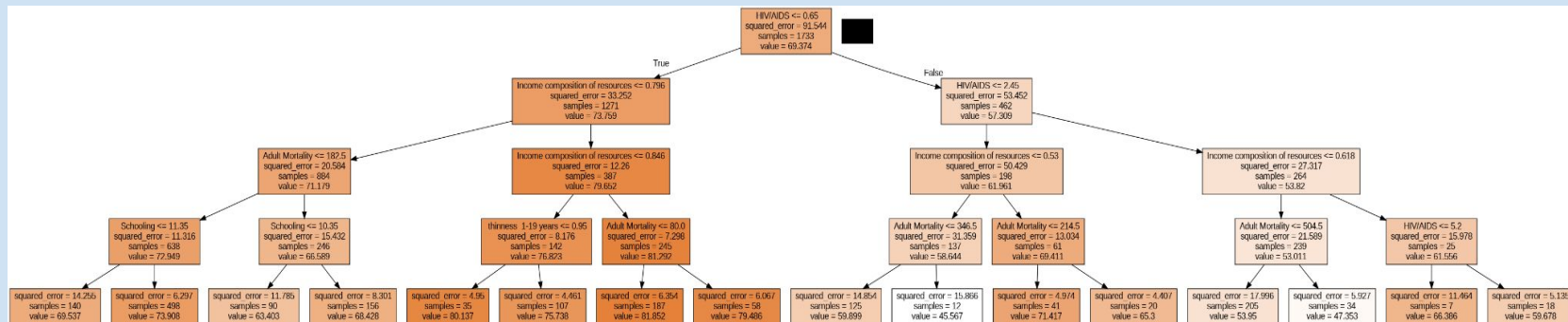
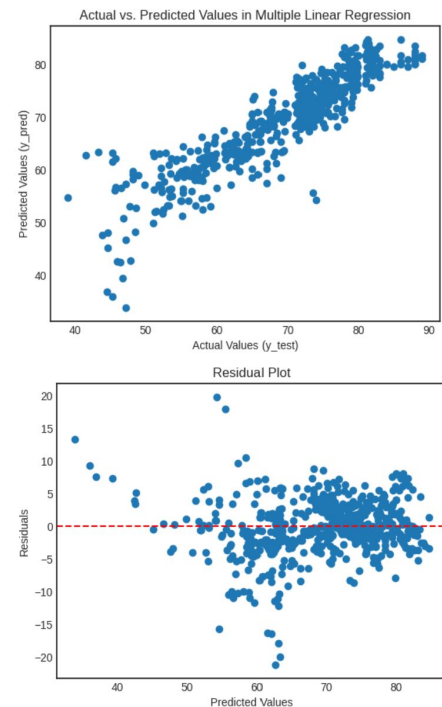
- Evaluation Metrics:

- Mean Absolute Error: 3.1173210283190875
    - Mean Squared Error: 18.682813144063324
    - Root Mean Squared Error: 4.322361986699324
    - R-squared: 0.8166738752230451

- Decision Tree

- Evaluation Metrics:

- Mean Absolute Error: 2.58380725154973
    - Mean Squared Error: 13.366842793673236
    - Root Mean Squared Error: 3.656069309199873
    - R-squared: 0.8688371247428789





# Comparison and Conclusion

- Overall, decision tree had better evaluation scores
- All metrics (RMSE, MAE, MSE,  $R^2$ ) were smaller in Decision Tree
- Important note:
  - Decision Tree will be restricted to the data within this range
  - To predict values outside of the range, MLR will be a better choice

# Model + Dataset

Model:

[https://colab.research.google.com/drive/1nUCQD8mjUMs-TKhQU4J8cs2tPx\\_UDEAA?usp=sharing](https://colab.research.google.com/drive/1nUCQD8mjUMs-TKhQU4J8cs2tPx_UDEAA?usp=sharing)

Dataset:

<https://drive.google.com/file/d/1cLYx7i3dMytWbhYU--g1TUilO9aklHHg/view?usp=sharing>