
LIFE EXPECTANCY

DOES IMMUNIZATION MATTER?

Mireille P. Feudjio T.
Springboard School of Data
Student

Outline

- **Context**
- **Data Wrangling**
- **Exploratory data analysis**
- **Modelling**
- **Conclusion**



Context

- Past 15 years, development of health section - reduction human mortality
 - Developing countries (30 years)
 - Assessment of population health – life expectancy
 - vaccination – improving life expectancy
 - pandemic
 - Impact of vaccination on life expectancy?
 - does immunization matter?
 - Factors contribution in life expectancy ?
-
- The dataset (life expectancy, health factors for 193 countries) 2000-2015.

within a month, the present project assessed the contribution of features on life expectancy with a special focus on immunization factors, and develop a regression model to predict life expectancy.

Constraints : The dataset has important missing values to handle .

Who care? Governments are concerned about the health of the population and are constantly looking for ways to improve life expectancy.

Data wrangling (1)

Dataset has 2938 observations and 22 columns (21 are independent variables).

Predicting variables were classified in 4 categories:

- Immunization related factors.
- Mortality factors.
- Economical factors.
- and Social factors.

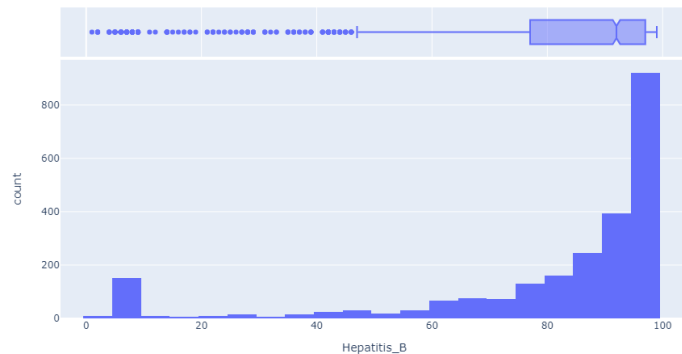
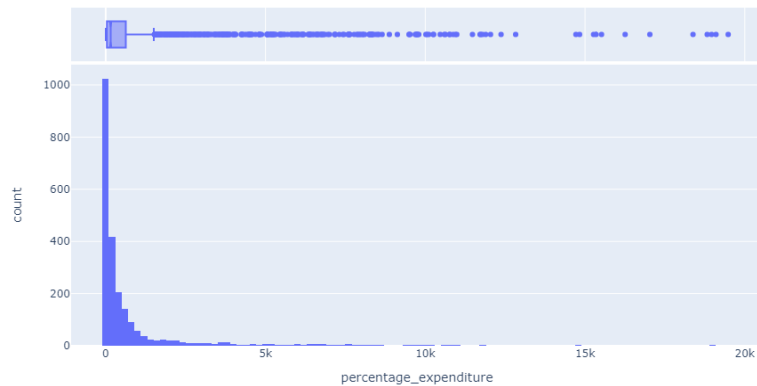
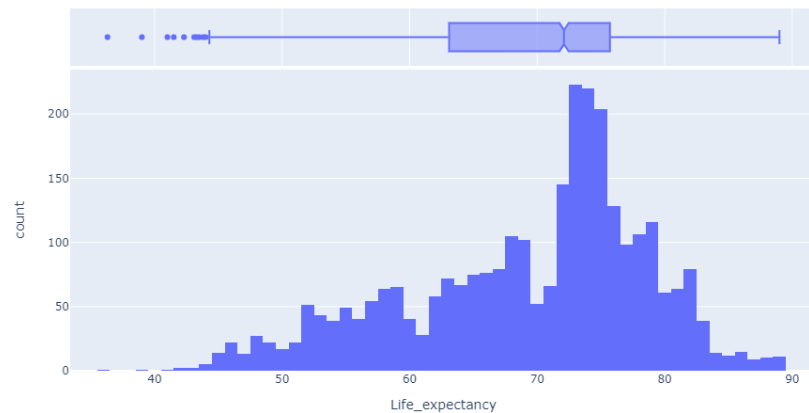
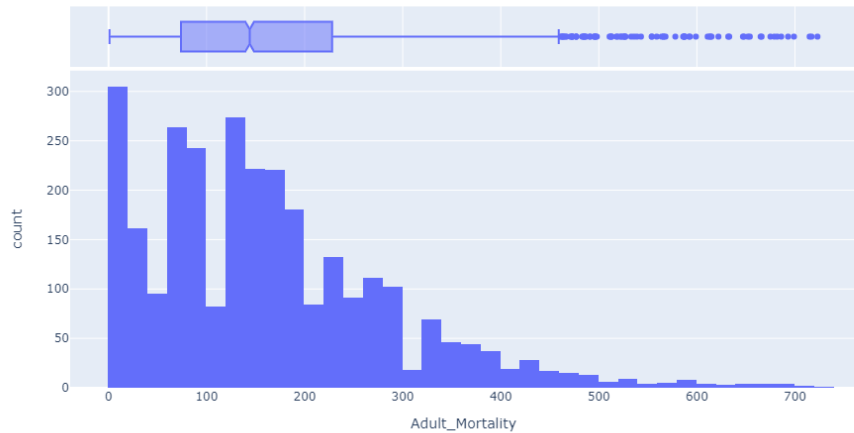
Data wrangling:

- Renaming column
- Replacing uncommon value with NAN
- Check for data type
- Check for duplicate
- Analysis missing value
- Analyzing outliers

<p>Mortality factors</p> <ul style="list-style-type: none">• Adult Mortality Rates of both sexes• Number of Infant Deaths• Number of under-five deaths• Deaths per 1 000 live births HIV/AIDS (0-4 years)• Measles - number of reported cases	<p>Social factors</p> <ul style="list-style-type: none">• Alcohol, recorded per capita (15+) consumption• Number of years of Schooling(years)• Average Body Mass Index• Prevalence of thinness -Age 10 to 19 (%))• Prevalence of thinness - Age 5 to 9(%)
<p>Immunization factors</p> <ul style="list-style-type: none">• Hepatitis B (HepB) immunization coverage• Polio (Pol3) immunization coverage• Diphtheria tetanus toxoid and pertussis (DTP3) immunization coverage	<p>Economic factors</p> <ul style="list-style-type: none">• Gross Domestic Product (in USD)• General government expenditure on health as a percentage of total government expenditure (%)• Expenditure on health as a percentage of Gross Domestic Product per capita (%)• Human Development Index in terms of income composition of resources (index ranging from 0 to 1)

Data wrangling (2)

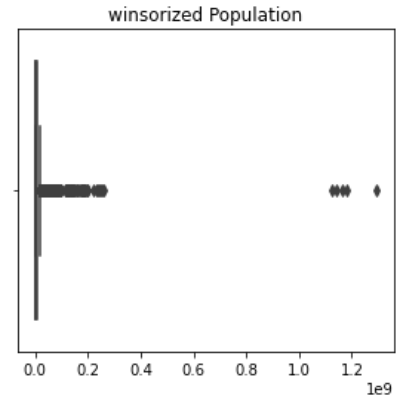
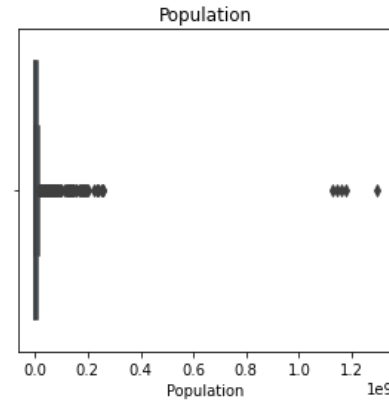
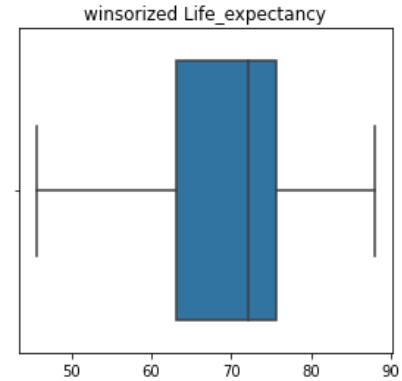
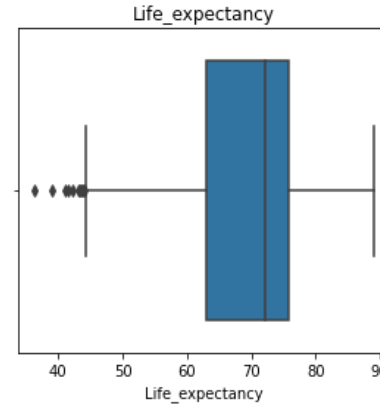
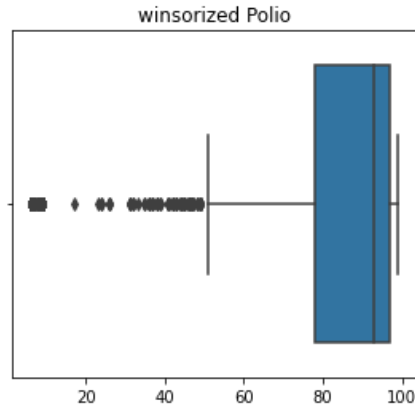
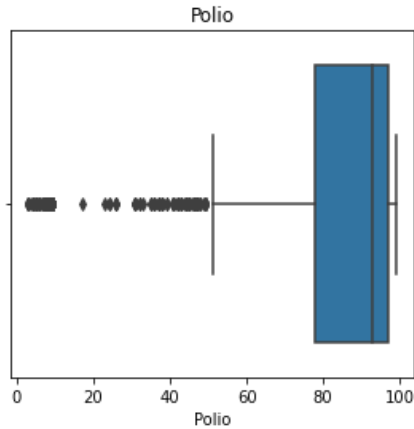
Outliers visualization and treatment



Data wrangling (3)

Outliers visualization and treatment

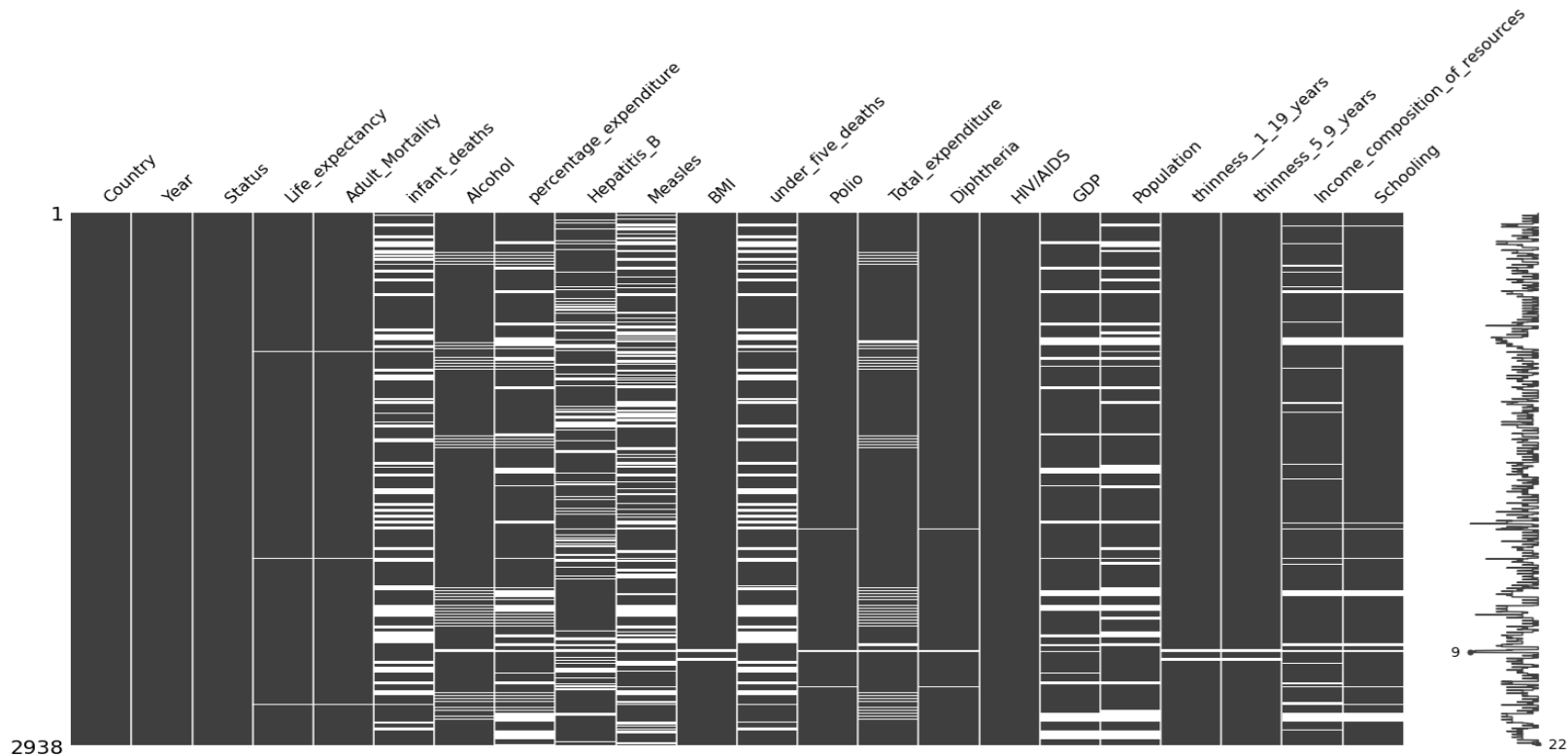
Winsorize method to treat Outliers



Box plot before and after the 98% winsorize application (1% top and 1% bottom)

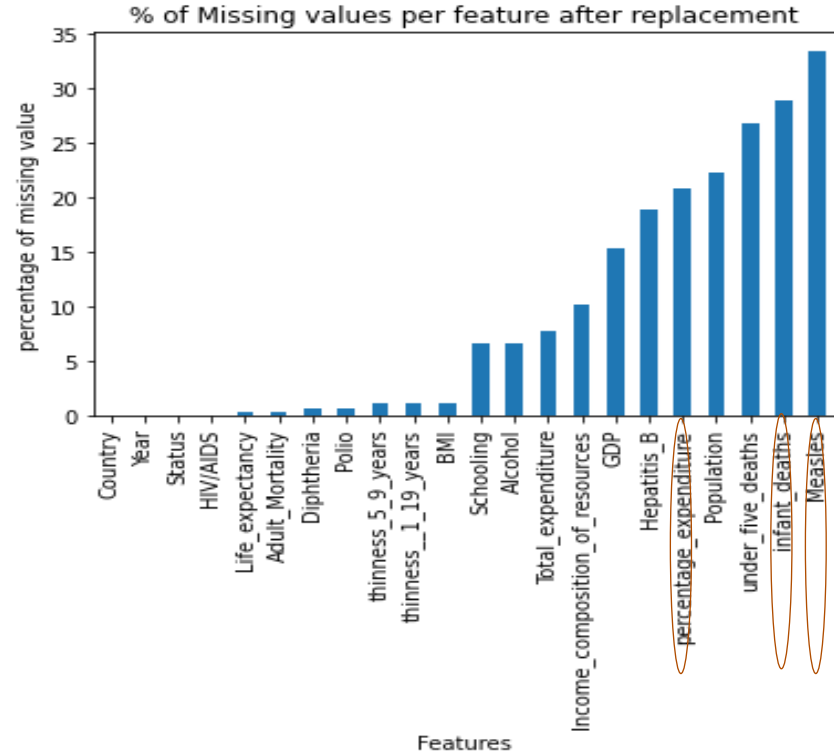
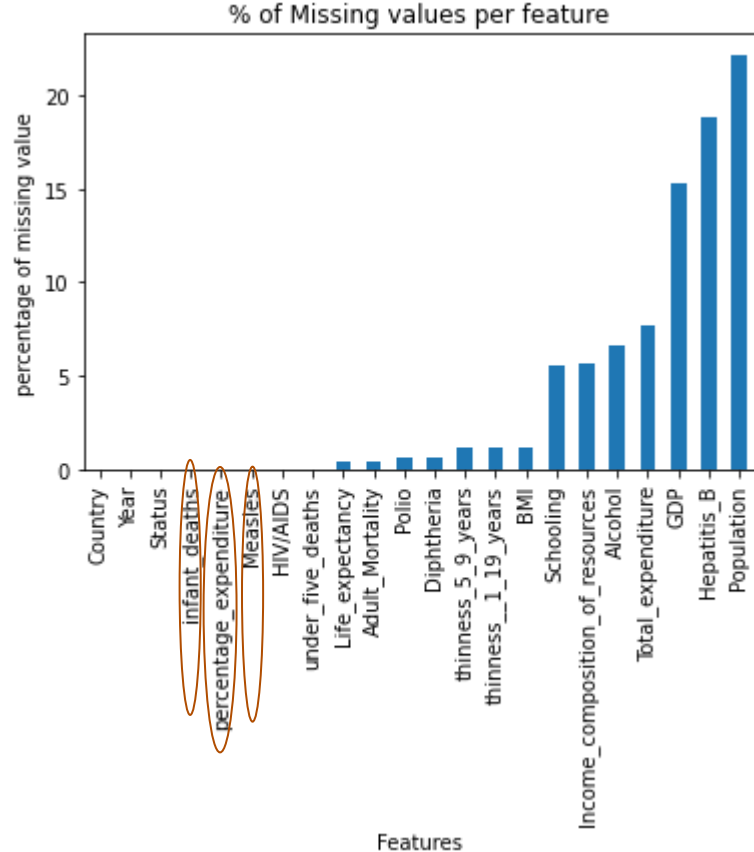
Data wrangling (4)

Assessing and treating missing value



Data wrangling (5)

Assessing and treating missing value



Original state of data with missing value and after the replacement of the uncommon type with NAN.

Data wrangling (6)

Assessing and treating missing value

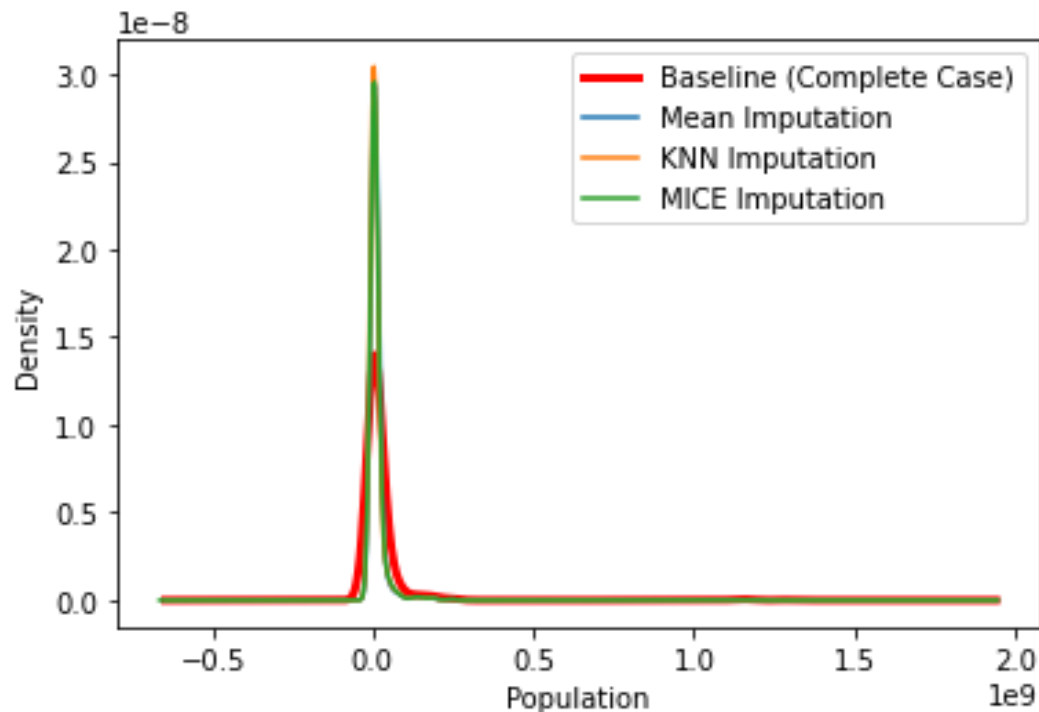
- loss of **69%**
- Mean, mode, constant and median imputation.
- KNN imputation and , MICE imputation using.
- Evaluation of the imputation methods.

Type of imputation	Adj. R-squared
Complete Case	0.899322
Mean Imputation	0.837261
KNN Imputation	0.857469
Mode Imputation	0.835953
Constant Imputation	0.665765
MICE Imputation	0.900985

Data wrangling (7)

Assessing and treating missing value

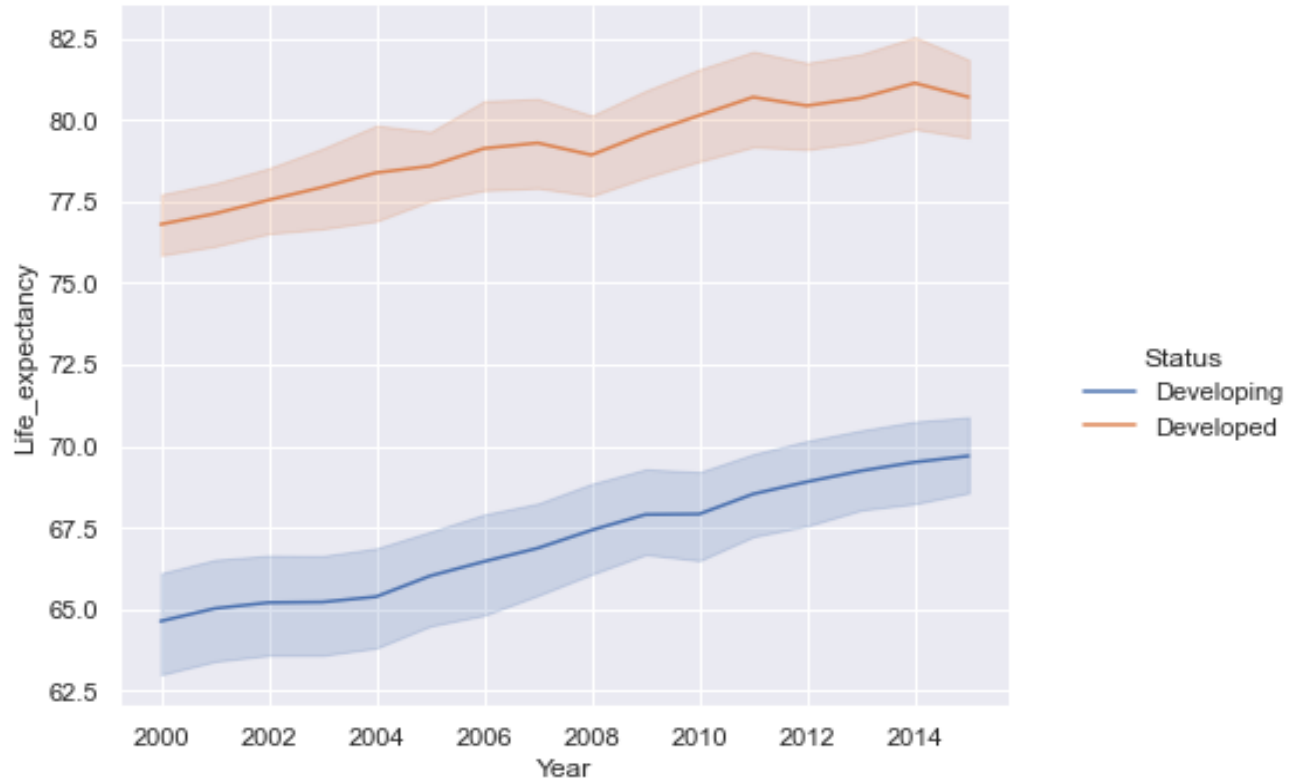
The best imputation technique is:
MICE Imputation



Density plot of population with three imputation methods

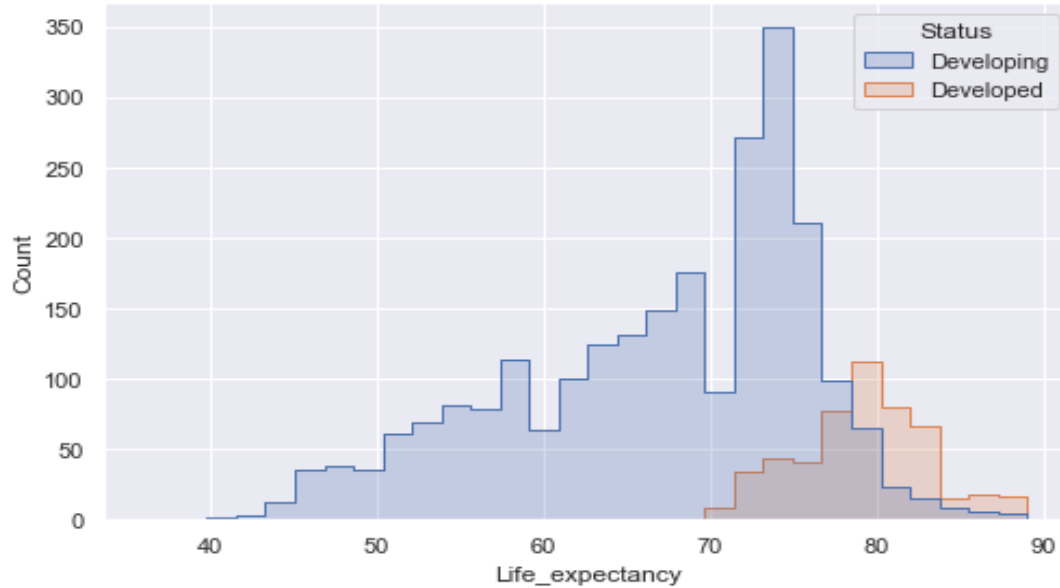
Exploratory Data Analysis (1)

What is the trend of life expectancy?



Exploratory Data Analysis (2)

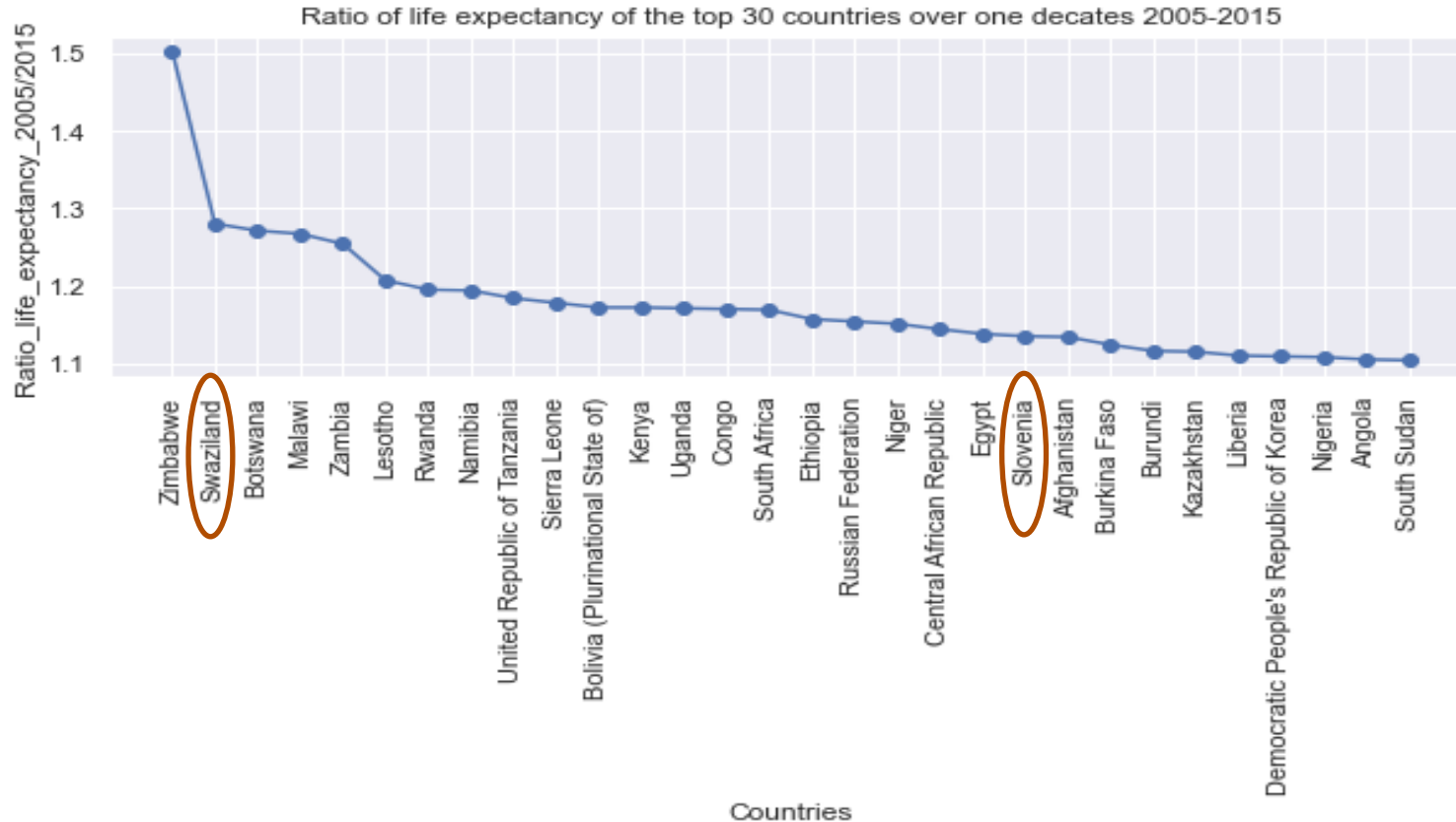
How is the distribution of life expectancy look like?



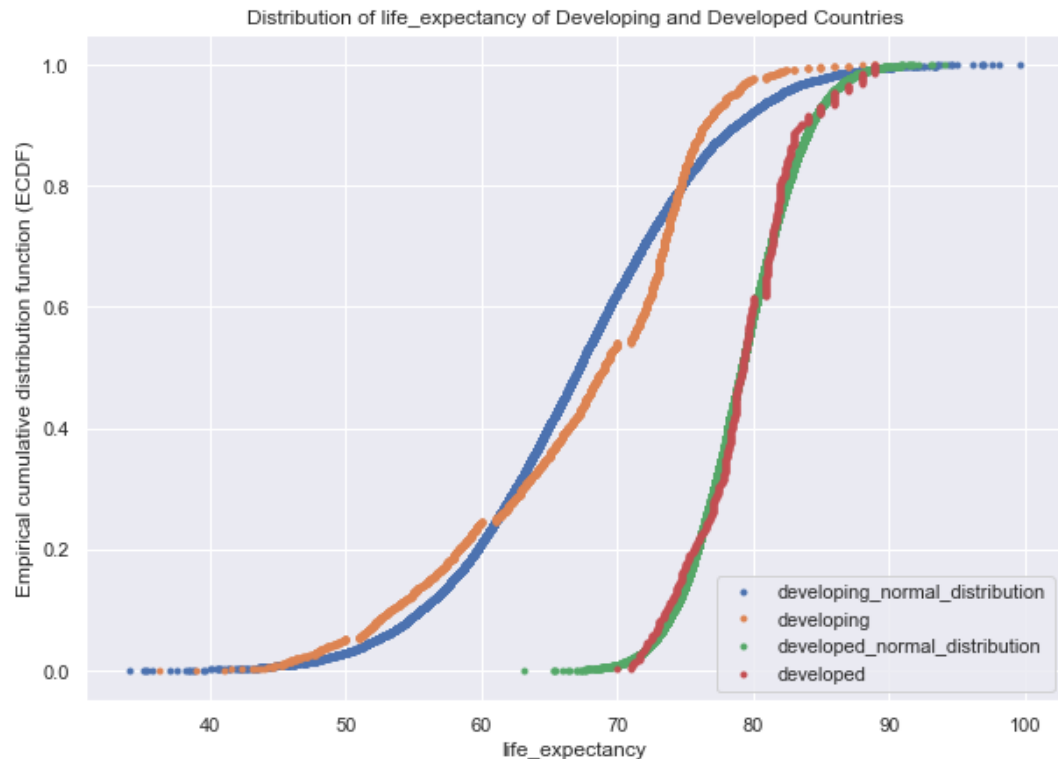
Countries types	count	mean	std	min	max
Developed	512.00	79.20	3.93	69.90	89.00
Developing	2416.00	67.11	9.00	36.30	89.00

Exploratory Data Analysis (3)

How was life expectancy over one decade (2005 to 2015) ?



What is the confidence interval of life expectancy at 95% ?



Apply the Central limit theorem

Confidence Interval of Life expectancy at 95%

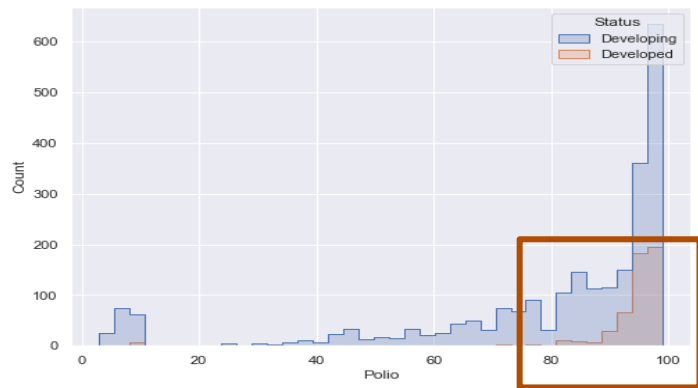
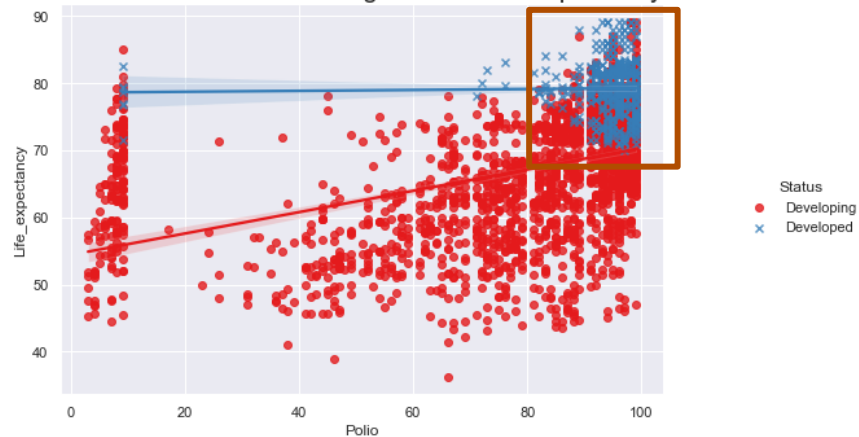
*CI developing countries : [64.834, 68.364]

*CI developed countries : [77.911, 79.451]

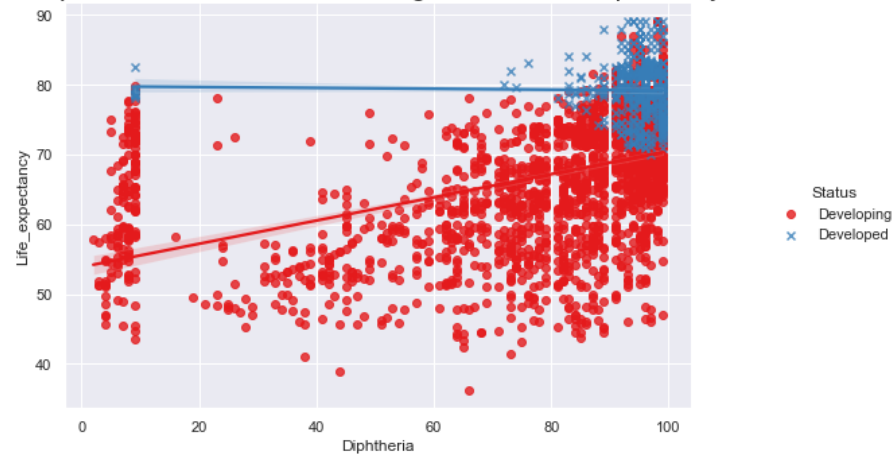
Exploratory Data Analysis (5)

Immunization and life expectancy

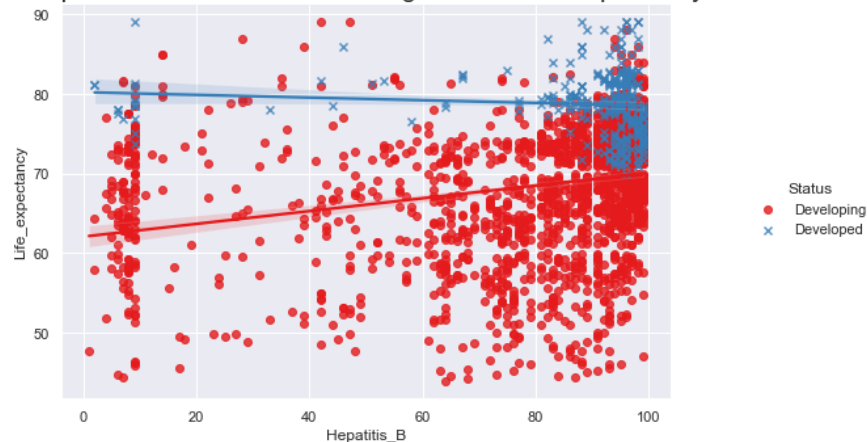
Polio immunization coverage versus Life expectancy



Diphtheria immunization coverage versus Life expectancy

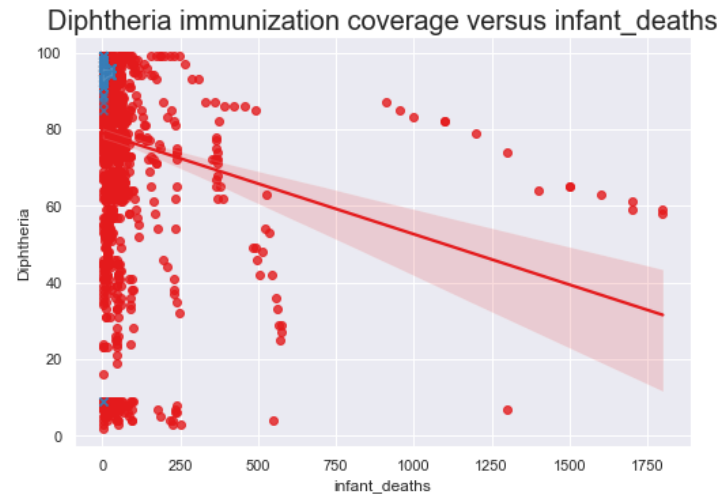
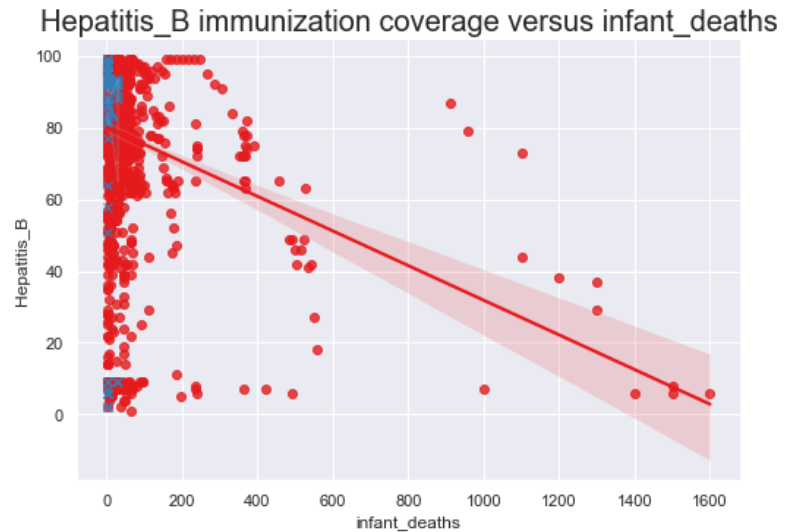
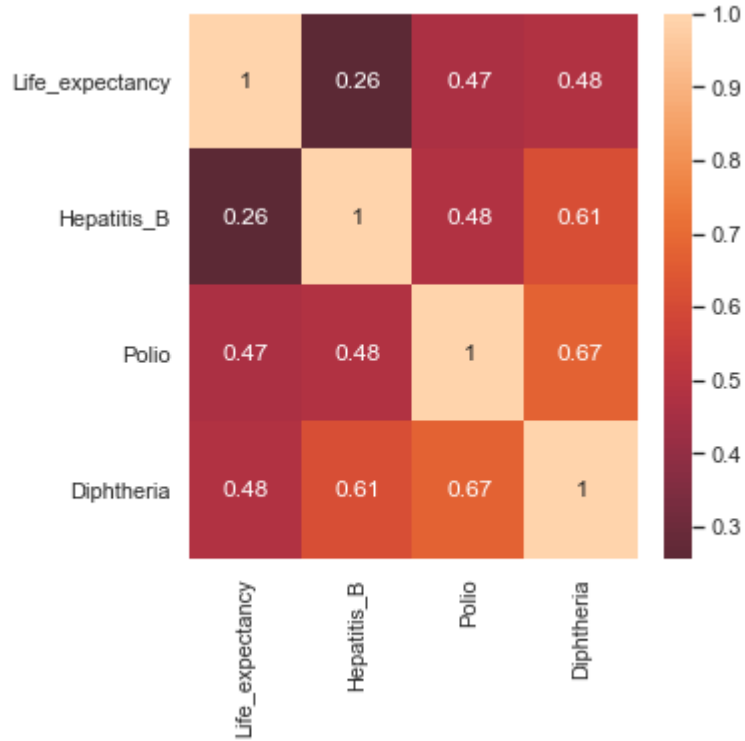


Hepatitis B immunization coverage versus Life expectancy



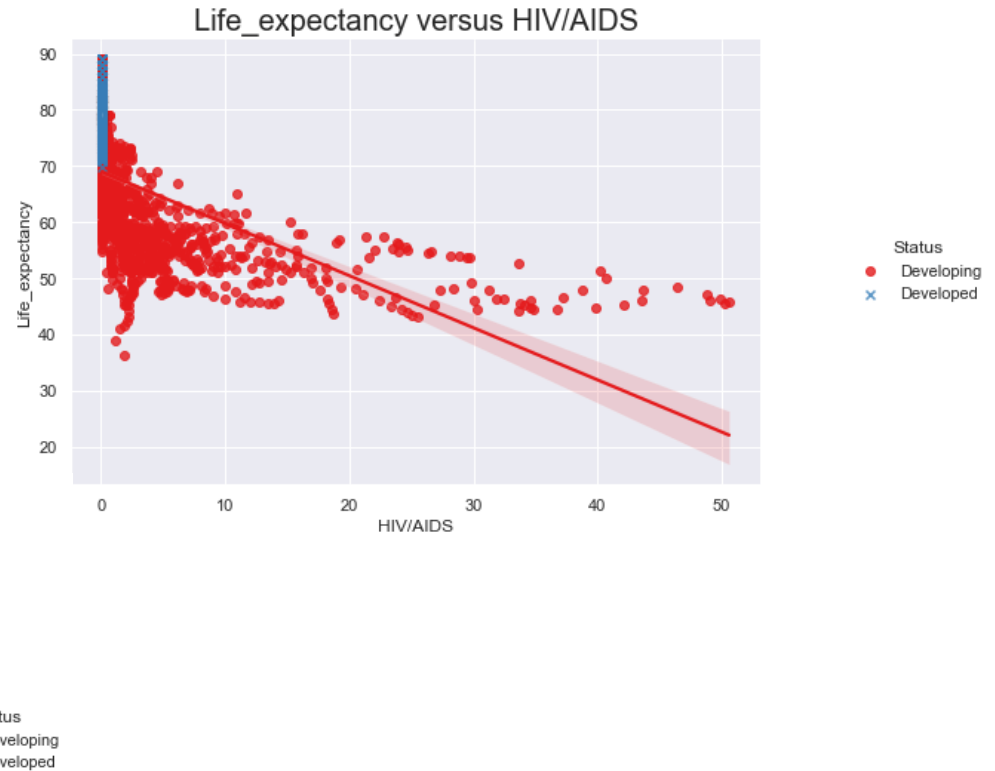
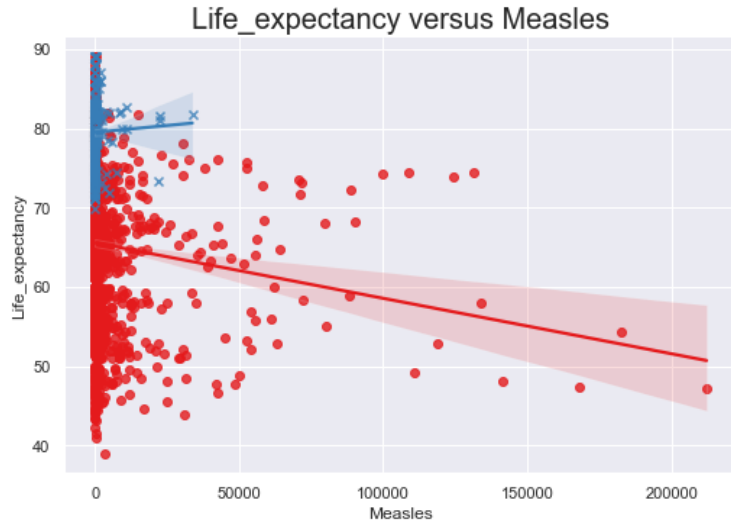
Exploratory Data Analysis (6)

Immunization and life expectancy , and mortality factors

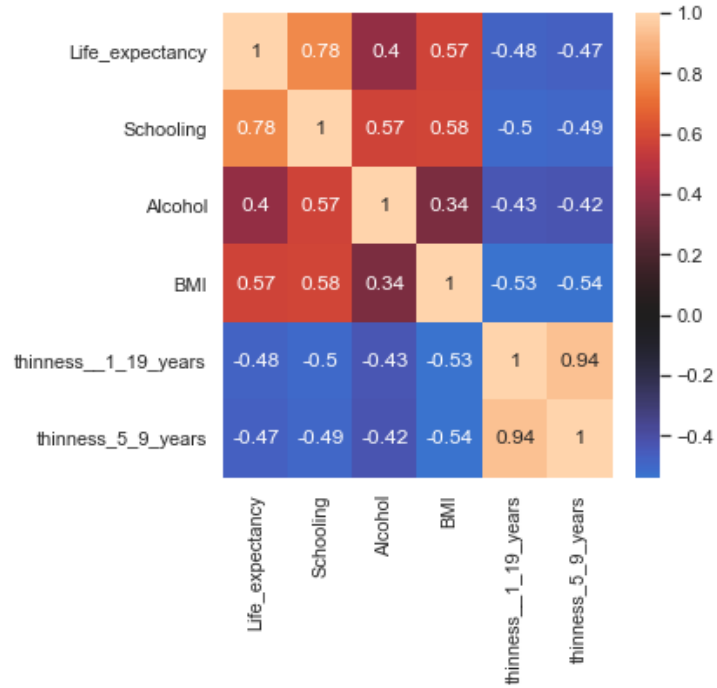


Exploratory Data Analysis (7)

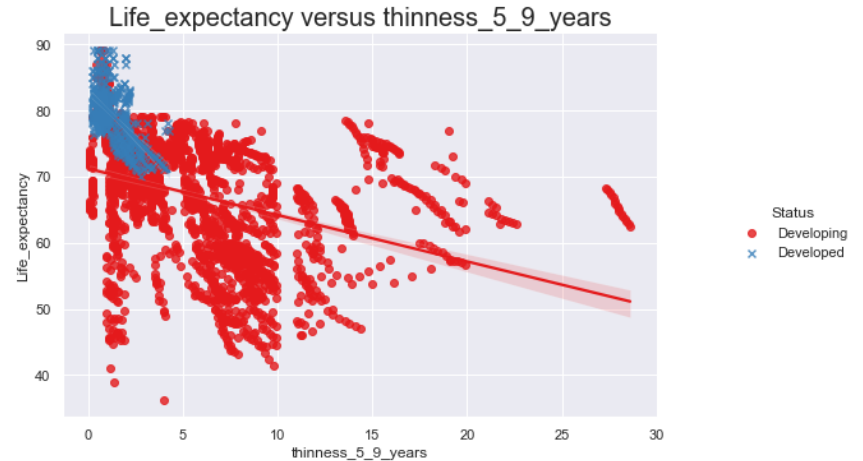
Mortality_factors and Life_expectancy



Exploratory Data Analysis (8)

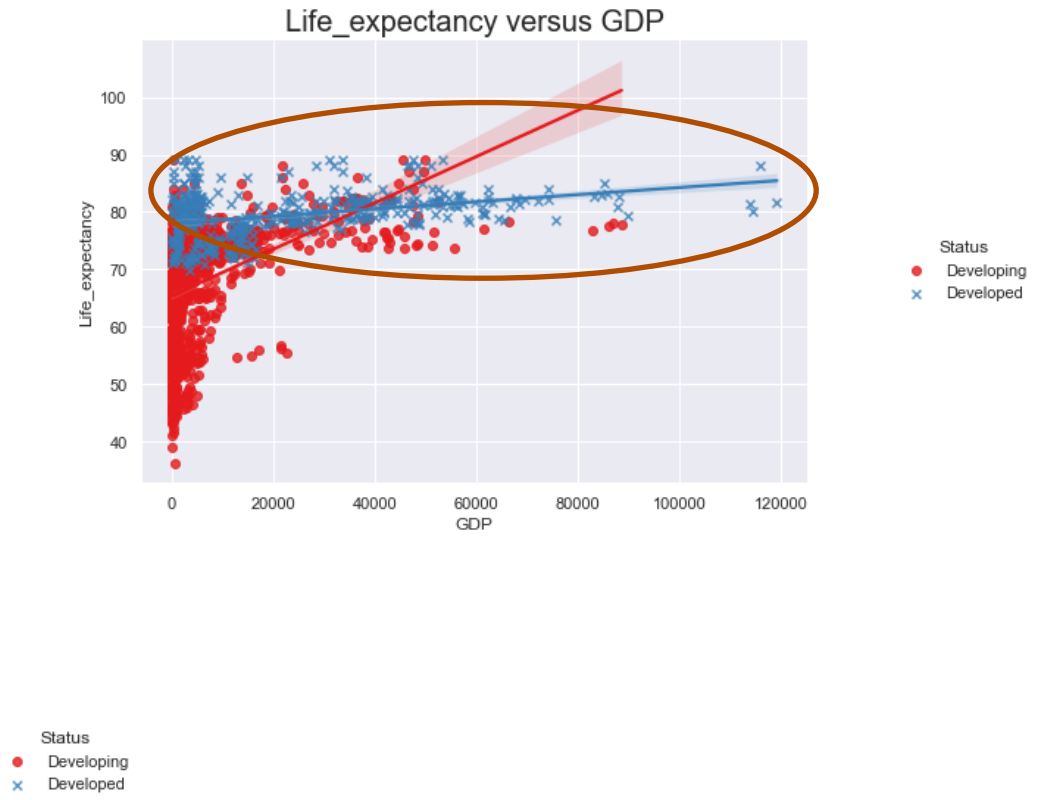
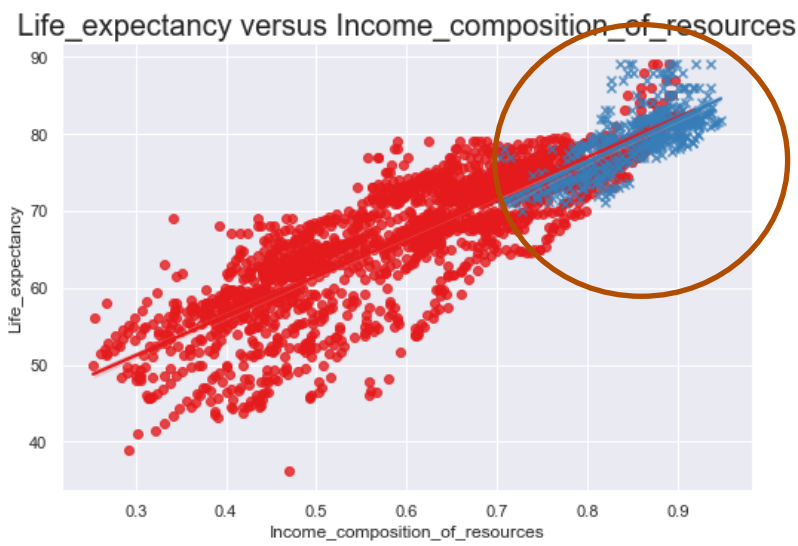


Social Factors and life_expectancy



Exploratory Data Analysis (9)

Economical_factors and Life_expectancy



Exploratory Data Analysis (10)

Population and life expectancy



Modeling (1)

- Label Encoder of categorical variable
- Split data into test set and train set (30% , 70%)
- Imputing missing value (Mice imputation)
- Scaling the dataset
- Model specifications
- GridSearchCV
- Random Search CV
- Train and evaluate the models, models selection

Modelling (2)

Linear Models

Model	Model definition
linear_reg	<pre>Pipeline(steps=[('iterativeimputer', IterativeImputer()), ('standardscaler', StandardScaler()), ('selectkbest', SelectKBest(k=22, score_func=<function f_regression at 0x000002477C893670>)), ('linearregression', LinearRegression())])</pre>
linear_reg2	<pre>Pipeline(steps=[('iterativeimputer', IterativeImputer()), ('standardscaler', StandardScaler()), ('pca', PCA(n_components=22)), ('linearregression', LinearRegression())])</pre>
ridge_reg	<pre>Pipeline(steps=[('iterativeimputer', IterativeImputer()), ('standardscaler', StandardScaler()), ('ridge', Ridge(alpha=0.5))])</pre>
Elastic_net	<pre>ElasticNet(alpha=0.0, l1_ratio=0.0)</pre>

Modelling (4)

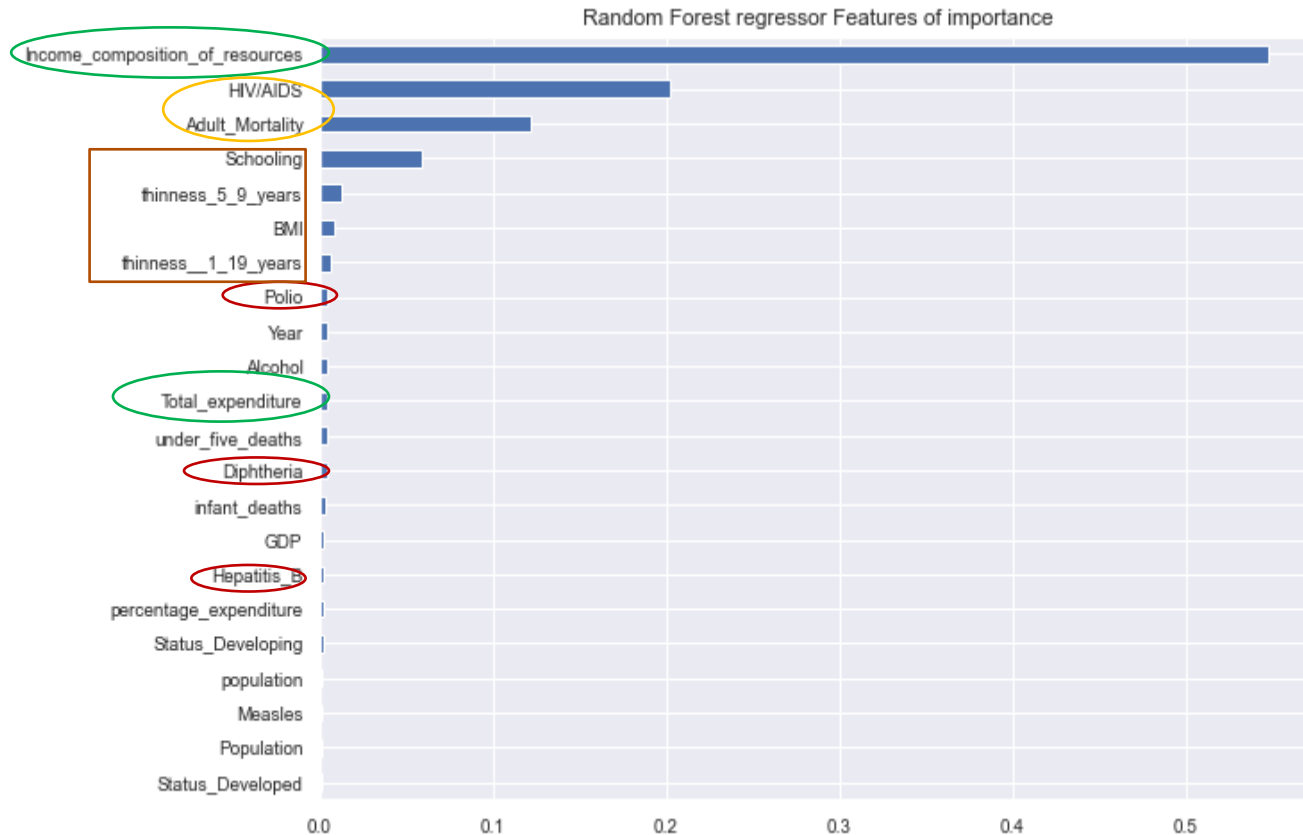


Modelling (5)

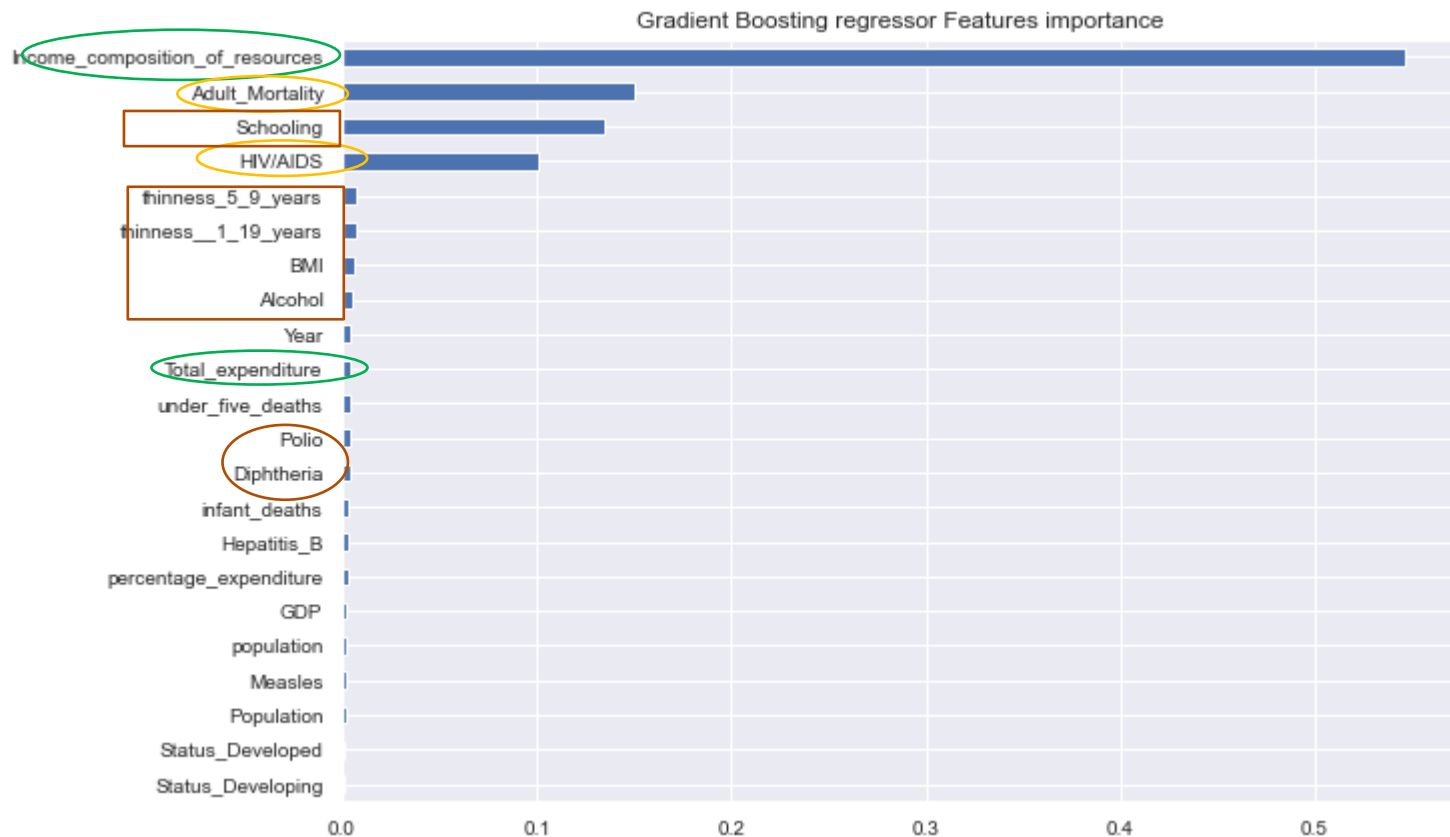
Tree based models

Model	Model definition
decision tree	DecisionTreeRegressor(max_depth=4, max_features=0.2, min_samples_leaf=0.1, random_state=1)
random_forest_reg1	RandomForestRegressor(max_depth=4, max_features=0.4, n_estimators=200, n_jobs=-1, random_state=1)
random_forest_reg2	RandomForestRegressor(max_depth=10, max_features=0.4, n_estimators=200, n_jobs=-1, random_state=1)
random_forest_reg3	RandomForestRegressor(max_depth=7, max_features=0.3, n_jobs=-1, random_state=1)
random_forest_reg4	RandomForestRegressor(max_depth=8, max_features=0.6, n_estimators=200, random_state=1)
gradien_boost_1	GradientBoostingRegressor(n_estimators=150, random_state=1)
gradien_boost_2	GradientBoostingRegressor(learning_rate=0.082, max_depth=10, max_features=0.60, min_samples_leaf=8, min_samples_split=10, n_estimators=118)
gradien_boost_3	GradientBoostingRegressor(learning_rate=0.082, max_depth=10, max_features=0.60, min_samples_leaf=8, min_samples_split=6, n_estimators=150, random_state=1)
voting	VotingRegressor(estimators=[('gb', GradientBoostingRegressor(random_state=47)), ('rf', RandomForestRegressor(random_state=47)), ('lr', LinearRegression())])
XGBRegressor	XGBRegressor(learning_rate=0.04, max_depth=5, n_estimators=200, n_jobs=-1, random_state=0)

Modelling (6)

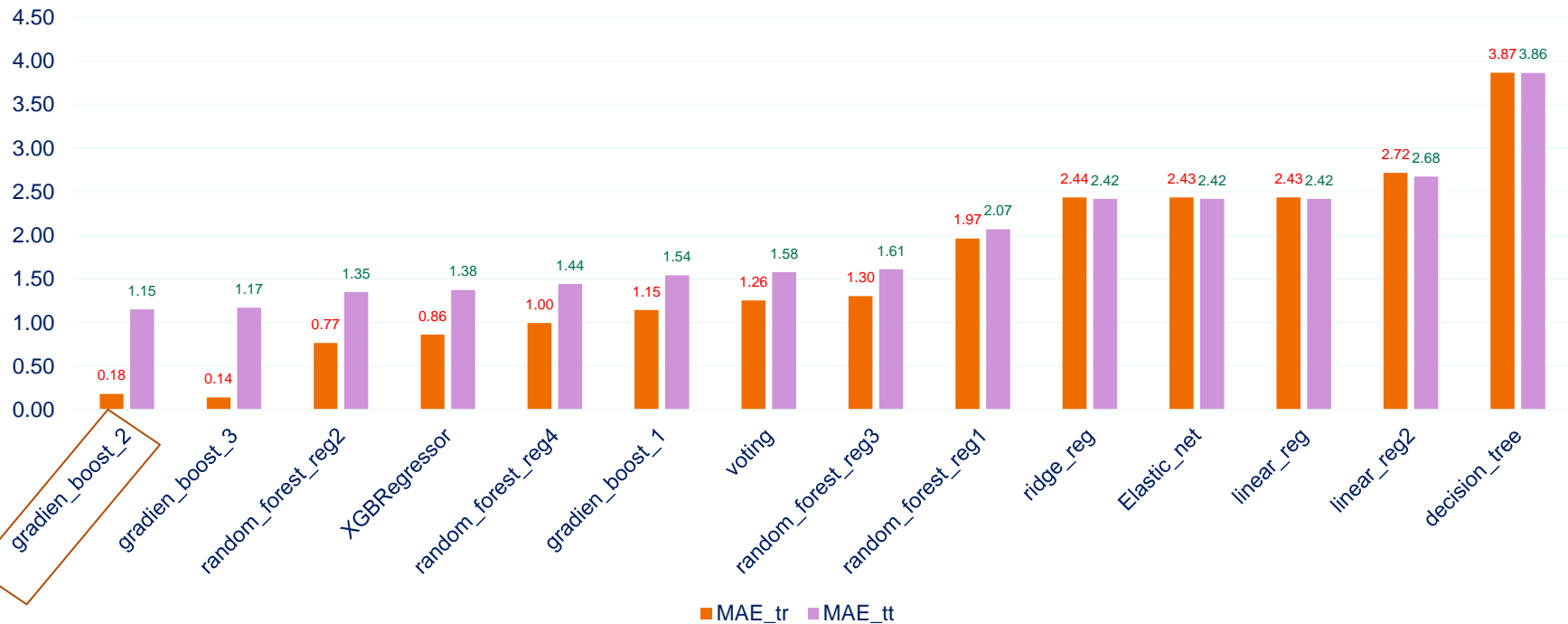


Modelling (7)



Modeling (8)

Model Evaluation- MAE of train/test set



Conclusion (1)

- Life expectancy has increased over years in both developed and developing countries.
- The mean average of the life expectancy of developed countries is generally higher compared to that of developing countries.
- However, the life expectancy ratio during the decade 2005-2015 showed that life expectancy in developing countries has increased significantly.
- It has been highlighted that immunization has impacted the improvement of life expectancy in a developing country, as well as the reduction in infant deaths.
- However, feature of importance reveals that immunization features has a very low contribution in the model .

Conclusion (2)

- The analysis revealed that economic factors and mortality factors play an important role in the system.
- It is why countries with higher income resources and GDP tend to have high life expectancy even if the population is big.
- In developing countries, an increase in the population tends to negatively impact life expectancy.
- (14) regression models have been developed to predict life expectancy.
- the chosen one is Gradient boost with **MAE of 0.18** on train set and **1.15 on the test** set.

Conclusion (3)

To increase life expectancy, the government must improve policies related to the human development index (income and resource composition) and social factors .

Although immunization factors are not the most important, they help reduce the effect of mortality factors on life expectancy.

Thus, promoting immunization will indirectly play a big role in life expectancy.

Thank you