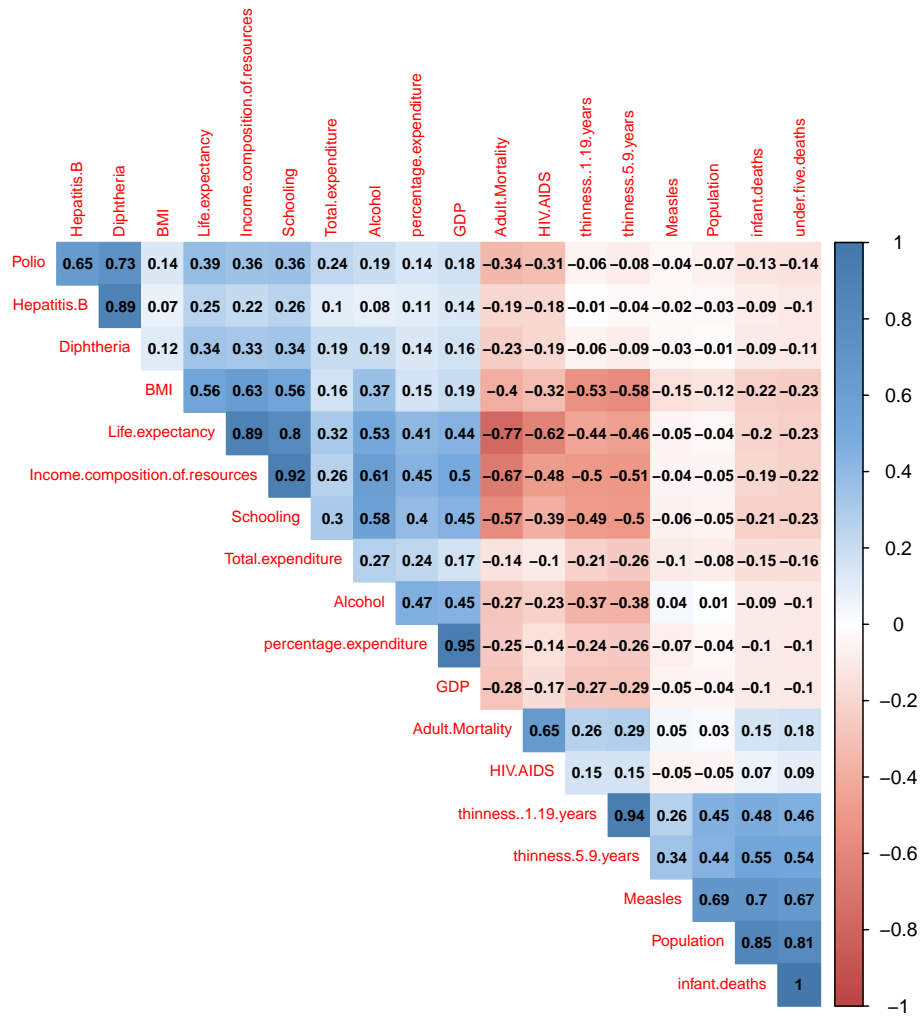# Primary relationship of interest

Team Orange

2022-10-19

## Primary relationship of interest
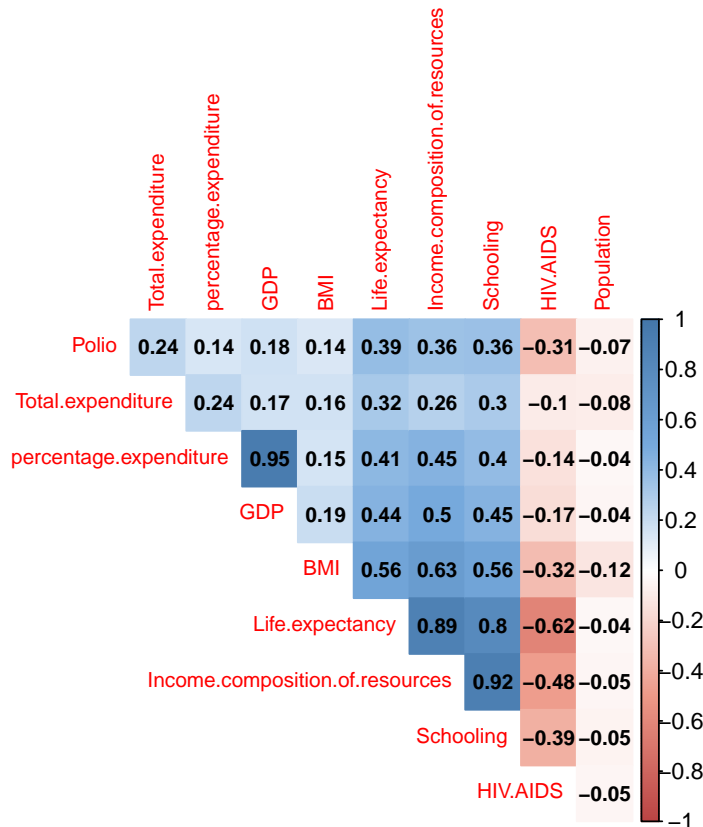
Full correlation map (Maybe in appendix):

**After piror selection:**

When building a model, it is not wise use all variables when they are highly correlated with each other. To represent immunization coverage, among "Hepatitis.B","Polio", "Diphtheria tetanus toxoid and pertussis (DTP3)", we decide to use "polio", since it has the highest correlation between Life Expectancy. Similarly, we see extremely high correlation between GDP and percentage expenditure, but both of them would have meaningful interpretation, we may want to decide which one goes to our final model when we conduct model selection.

On the one hand, with domain knowledge, we know "Adult.Mortality", "infant.deaths" and "under.five.deaths" variables are directly correlated to Life Expectancy, we choose to drop them from the predictor variable list; on the other hand, we are interested in "HIV.AIDS" variable (Deaths per 1 000 live births HIV/AIDS (0-4 years), .... (please add some reasons here)

We would like to omit the variables have Low correlation between Life Expectancy: "Measles"; however, we do want to include "population" because of our interest.

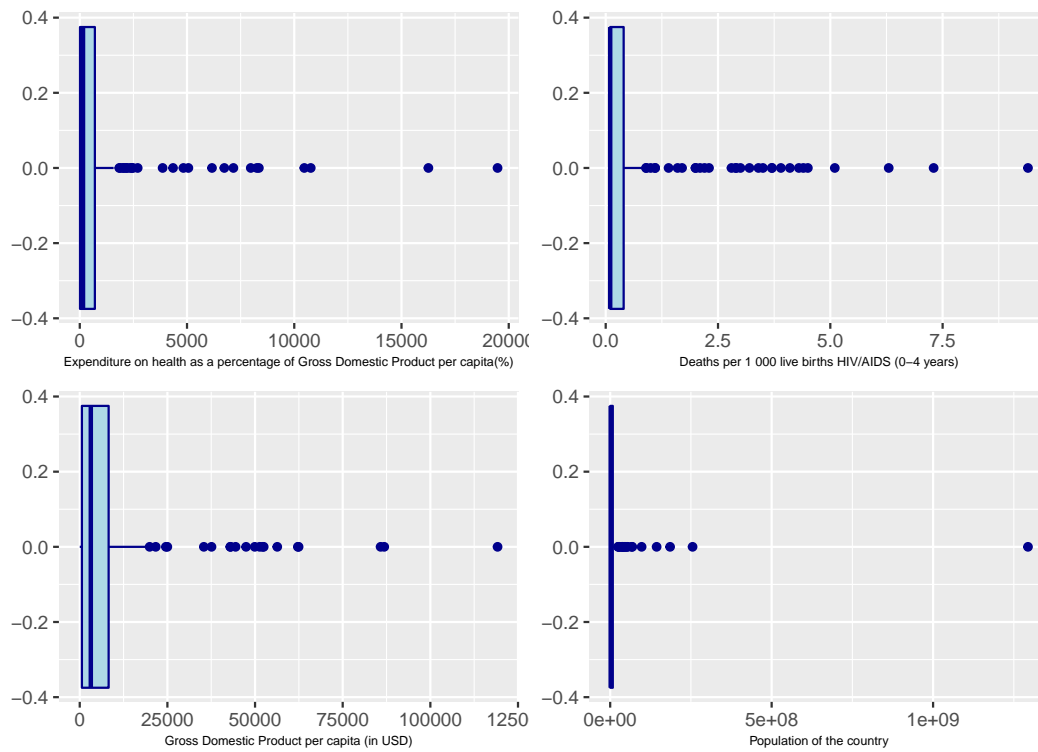As for categorical variable, we would like to keep country status (developing/developed) as one of the predictors



**Potential predictor variables:**

- *Percentage Expenditure* or *GDP*
- *BMI*
- *Polio*
- *HIV/AIDS*
- *Total.expenditure*

- *Schooling*
- *Income Composition of Resources*
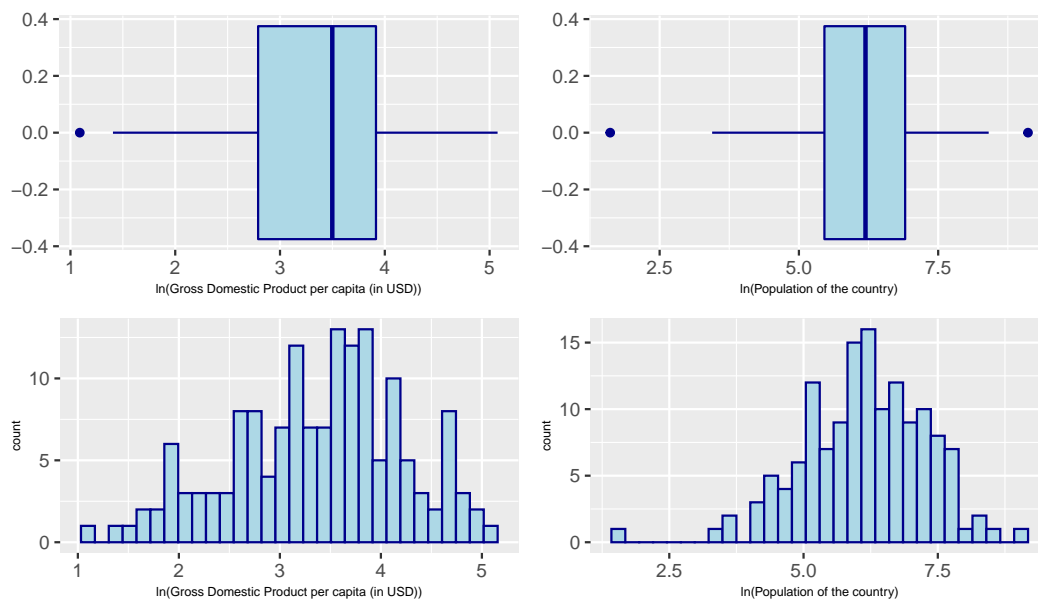- *Country Status (Developed/Developing)*
- *Population*

**Transformation if needed when modeling**

When checking box plots of all variables, we find 4 variables may need some transformation.



We may want to use log transformation for population and GDP, since the magnitude of gaps are huge. But for other two, we need more investigation due to the difficulty of interpretation.

The box plots and histograms of population and GDP after log transformation:
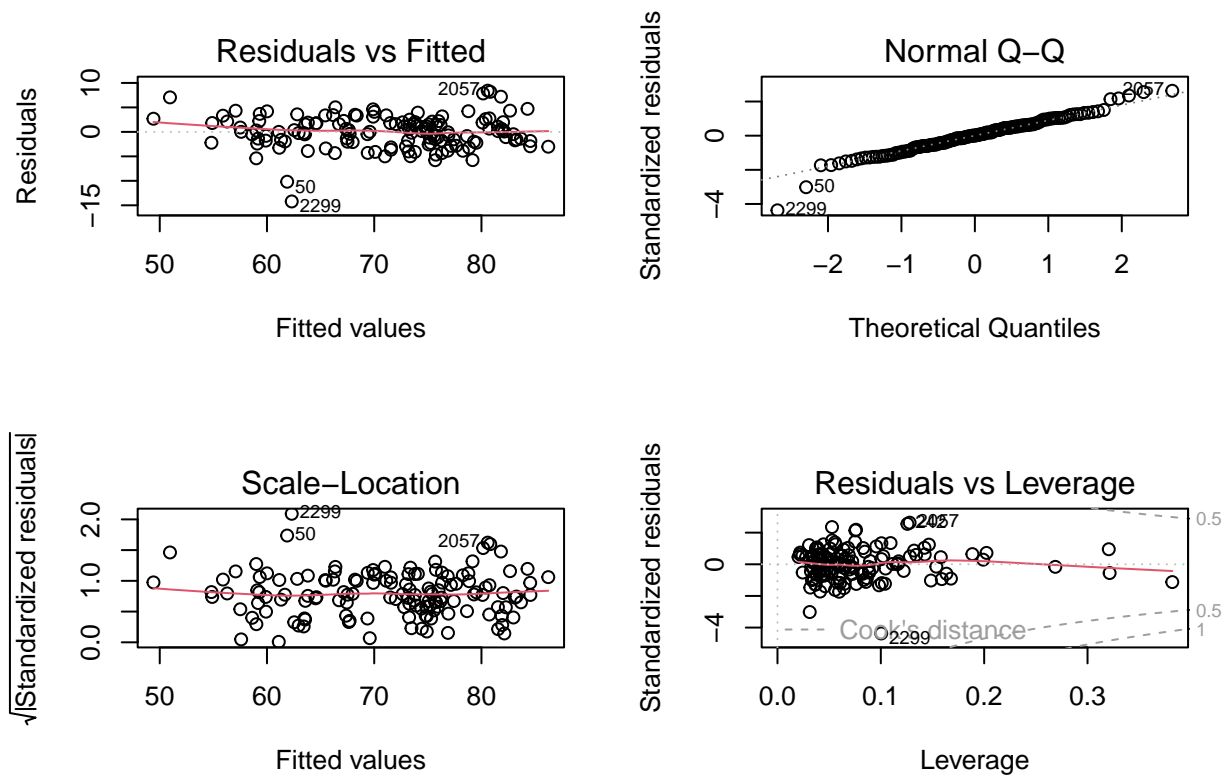
(For our reference)

After selection Model:

Table 1: Regression Summary

|  | Dependent variable: |
| --- | --- |
|  | Life.expectancy |
| BMI | −0.002 (0.017) |
|  | p = 0.925 |
| log10(GDP) | −0.219 (0.521) |
|  | p = 0.675 |
| percentage.expenditure | 0.0001 (0.0001) |
|  | p = 0.463 |
| Polio | 0.009 (0.015) |
|  | p = 0.557 |
| HIV.AIDS | −1.339 (0.228) |
|  | p = 0.00000$^{***}$ |
| Total.expenditure | 0.273 (0.120) |
|  | p = 0.025$^{**}$ |
| log10(Population) | −0.293 (0.262) |
|  | p = 0.267 |
| Income.composition.of.resources | 43.377 (5.893) |
|  | p = 0.000$^{***}$ |
| StatusDeveloping | −0.560 (1.010) |
|  | p = 0.580 |
| Schooling | −0.069 (0.271) |
|  | p = 0.800 |
| Constant | 44.058 (3.245) |
|  | p = 0.000$^{***}$ |
| Observations | 139 |
| $R^2$ | 0.863 |
| Adjusted $R^2$ | 0.852 |
| Residual Std. Error | 3.428 (df = 128) |
| F Statistic | 80.544$^{***}$ (df = 10; 128) |

| *Note:* | $^{*}$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01 |
| --- | --- |

Potential Final Model :

Table 2: Regression Summary

|  | *Dependent variable:* |
| --- | --- |
|  | Life.expectancy |
| log10(GDP) | −0.201 (0.422) |
|  | p = 0.634 |
| HIV.AIDS | −1.372 (0.214) |
|  | p = 0.000*** |
| Income.composition.of.resources | 42.946 (2.947) |
|  | p = 0.000*** |
| StatusDeveloping | −1.259 (0.853) |
|  | p = 0.143 |
| Constant | 44.597 (2.177) |
|  | p = 0.000*** |
| Observations | 154 |
| $R^2$ | 0.860 |
| Adjusted $R^2$ | 0.856 |
| Residual Std. Error | 3.326 (df = 149) |
| F Statistic | 228.657*** (df = 4; 149) |
| *Note:* | *p<0.1; **p<0.05; ***p<0.01 |