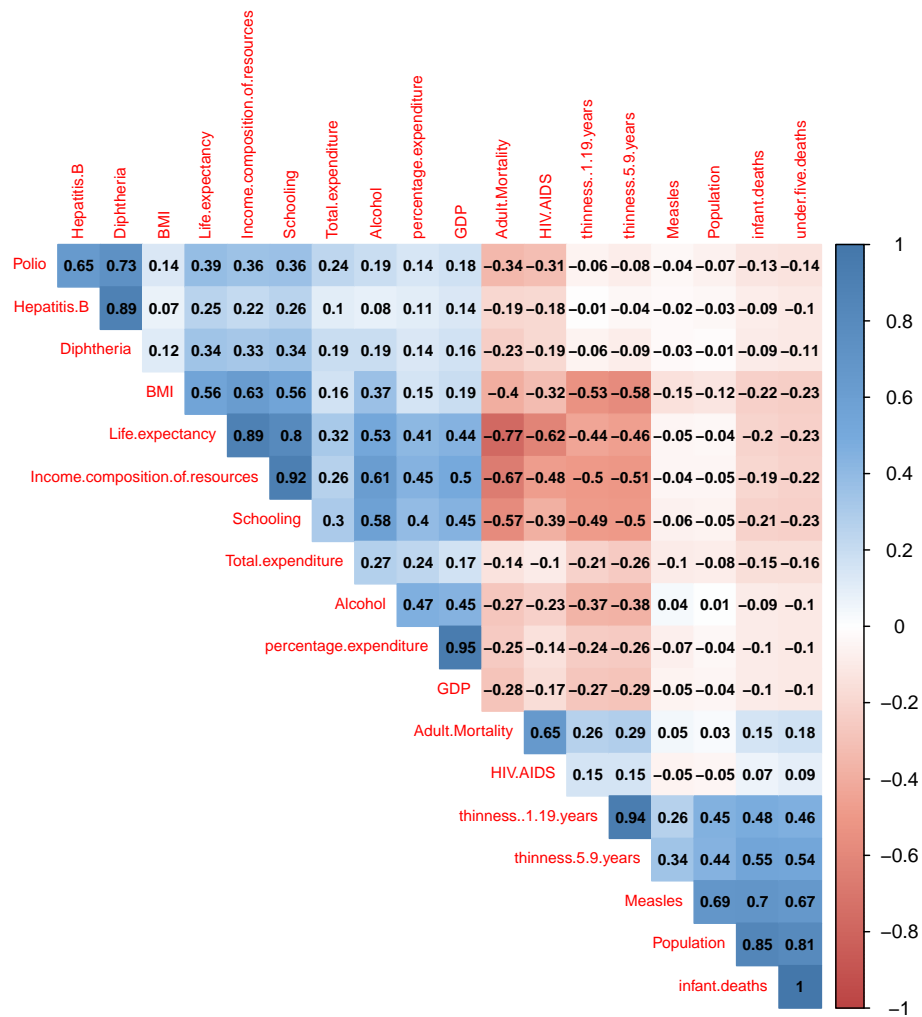# Primary relationship of interest

Team Orange

2022-10-19

## Primary relationship of interest

Full correlation map (Maybe in appendix):
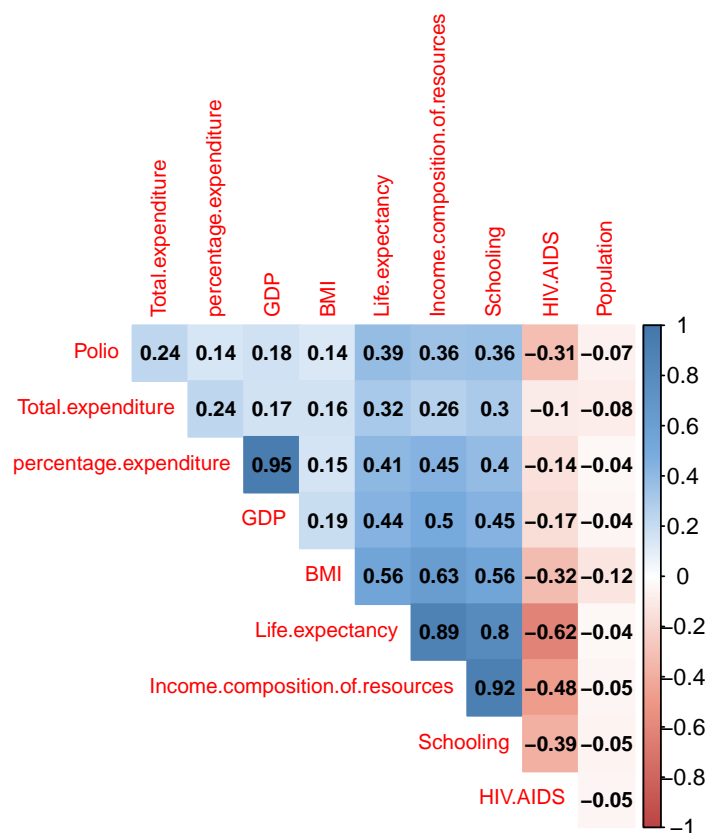
**After piror selection:**

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 interpreation, 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)

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.

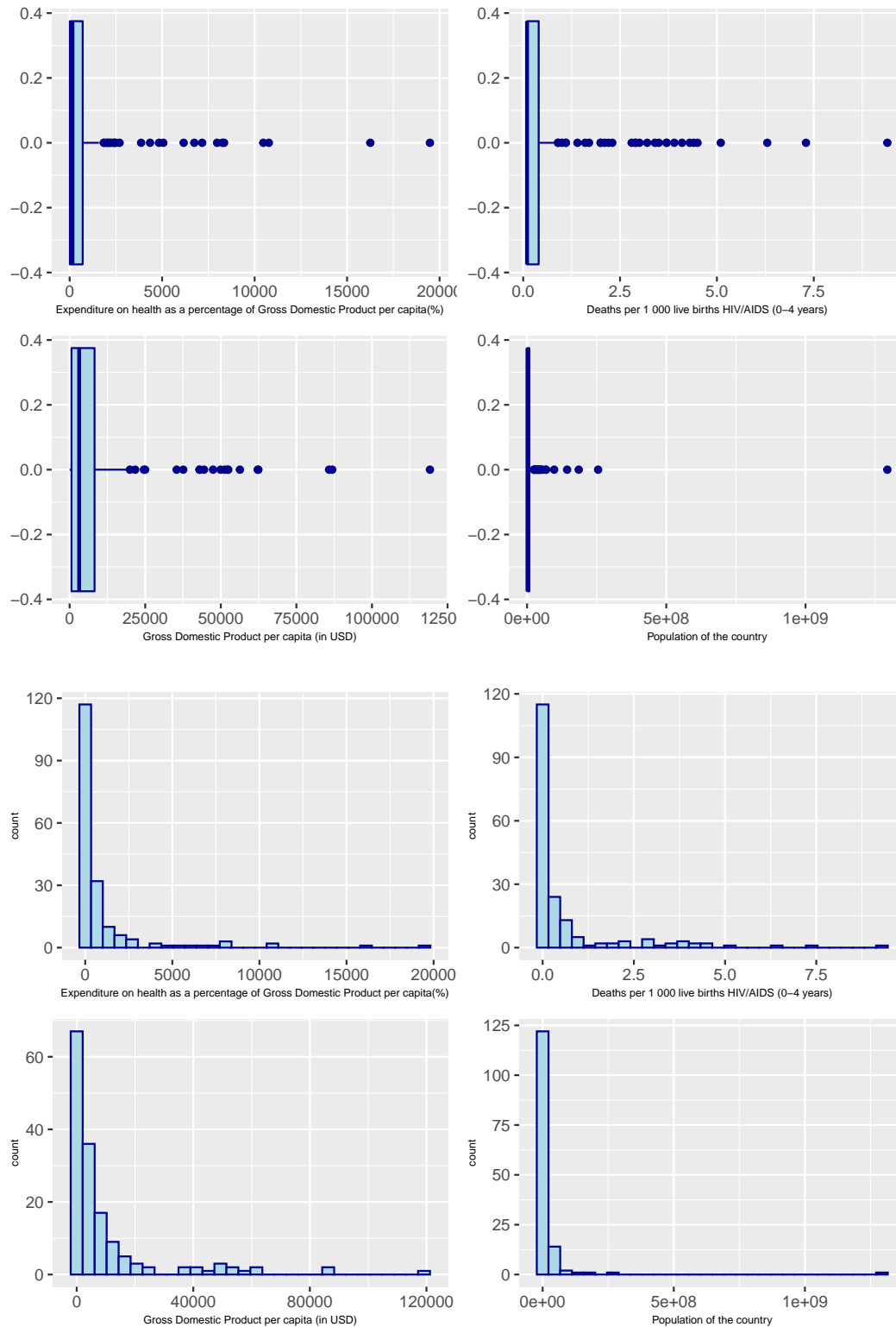As for categorical variable, we would like to keep country status (developing/developed) as one of the predictors



**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

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, plus the difficulty of interpretation.
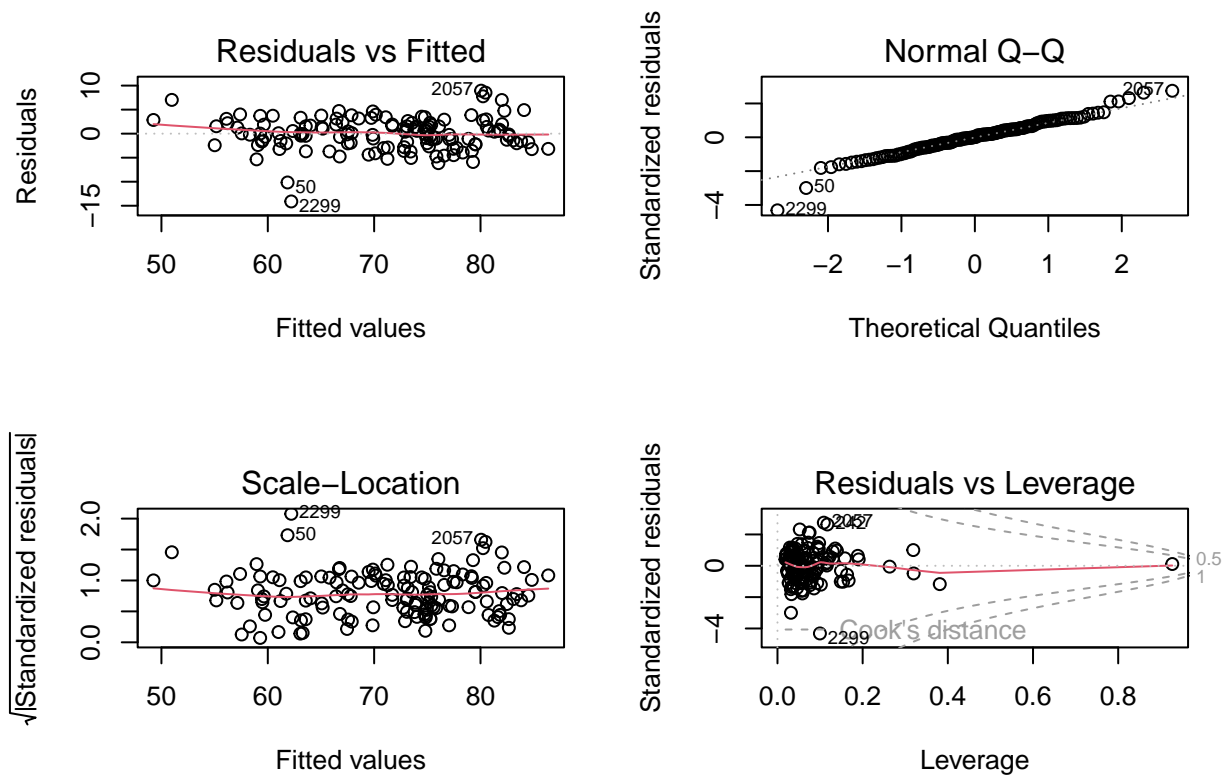
(For our reference)

After selection Model:

Table 1: Regression Summary

|  | Dependent variable: |
| --- | --- |
|  | Life.expectancy |
| BMI | $-0.001$ (0.018) |
|  | p = 0.942 |
| log10(GDP) | $-0.215$ (0.525) |
|  | p = 0.683 |
| percentage.expenditure | 0.0001 (0.0001) |
|  | p = 0.534 |
| Polio | 0.009 (0.015) |
|  | p = 0.535 |
| HIV.AIDS | $-1.347$ (0.230) |
|  | p = 0.00000*** |
| Total.expenditure | 0.282 (0.120) |
|  | p = 0.021** |
| Population | $-0.000$ (0.000) |
|  | p = 0.920 |
| Income.composition.of.resources | 44.229 (5.872) |
|  | p = 0.000*** |
| StatusDeveloping | $-0.574$ (1.015) |
|  | p = 0.573 |
| Schooling | $-0.106$ (0.270) |
|  | p = 0.697 |
| Constant | 42.063 (2.718) |
|  | p = 0.000*** |
| Observations | 139 |
| $R^2$ | 0.862 |
| Adjusted $R^2$ | 0.851 |
| Residual Std. Error | 3.444 (df = 128) |
| F Statistic | 79.651*** (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 |