

Final project: Step 1

Danyu Zhang, Limingrui Wan, Daniel Alonso

December 9th, 2020

Importing libraries

```
library(dplyr)
library(ggplot2)
library(reshape2)
library(PerformanceAnalytics)
library(gridExtra)
library(stringr)
library(foreach)
library(MASS)
library(andrews)
library(mice)
```

Importing data

```
data <- read.csv('./data/data.csv')
head(data)
#>   X      continent      location total_cases new_cases new_cases_smoothed
#> 1 0         Asia      Afghanistan    41728      95          99.429
#> 2 1         Africa      Angola      11035     230         236.286
#> 3 2         Europe      Albania     21523     321         296.857
#> 4 3         Europe      Andorra       4888      63          80.429
#> 5 4         Asia United Arab Emirates 135141    1234        1272.429
#> 6 5 South America      Argentina 1183118    9598        11547.143
#>   total_deaths new_deaths new_deaths_smoothed total_cases_per_million
#> 1         1544         3           3.143          1071.918
#> 2          286         2           2.571           335.755
#> 3          527         9           6.714          7478.977
#> 4           75         0           0.429         63262.797
#> 5          497         1           2.429        13663.856
#> 6        31623        483          331.714        26177.623
#>   new_cases_per_million new_cases_smoothed_per_million total_deaths_per_million
#> 1             2.440              2.554              39.663
#> 2             6.998              7.189              8.702
#> 3            111.544             103.154             183.126
#> 4            815.376            1040.944             970.685
#> 5            124.767             128.653              50.251
#> 6            212.365             255.492             699.689
#>   new_deaths_per_million stringency_index population population_density
#> 1             0.077           5.56    38928341           54.422
#> 2             0.061           NA    32866268           23.890
#> 3             3.127          50.93    2877800           104.871
#> 4             0.000          59.26     77265           163.755
```

```

#> 5          0.101          47.22    9890400          112.442
#> 6          10.687          81.94   45195777          16.177
#>   median_age aged_65_older aged_70_older gdp_per_capita extreme_poverty
#> 1         18.6         2.581         1.337         1803.987          NA
#> 2         16.8         2.405         1.362         5819.495          NA
#> 3         38.0        13.188         8.643        11803.431          1.1
#> 4          NA          NA          NA          NA          NA
#> 5         34.0         1.144         0.526        67293.483          NA
#> 6         31.9        11.198         7.441        18933.907          0.6
#>   cardiovasc_death_rate diabetes_prevalence hospital_beds_per_thousand
#> 1          597.029           9.59           0.50
#> 2          276.045           3.94           NA
#> 3          304.195          10.08           2.89
#> 4          109.135           7.97           NA
#> 5          317.840          17.26           1.20
#> 6          191.032           5.50           5.00
#>   life_expectancy human_development_index development
#> 1          64.83           0.498          low
#> 2          61.15           0.581         medium
#> 3          78.57           0.785          high
#> 4          83.73           0.858        very high
#> 5          77.97           0.863        very high
#> 6          76.67           0.825        very high

```

Excluding smoothed columns as they are redundant transformations of other columns

```

removed_cols <- c('new_deaths_smoothed', 'new_cases_smoothed', 'new_cases_smoothed_per_million', 'total_cases_smoothed')
data_n <- data
for (col in removed_cols) {data_n <- data_n[names(data_n) != col]}

```

Exploratory data analysis

Variable types

Categorical variables

- continent
- location
- development

Numerical variables

Discrete

- total_cases
- new_cases
- total_deaths
- new_deaths
- population

Continuous

- new_cases_smoothed
- new_deaths_smoothed
- total_cases_per_million
- new_cases_per_million
- new_cases_smoothed_per_million
- total_deaths_per_million
- new_deaths_per_million
- stringency_index
- population_density
- median_age
- aged_65_older
- aged_70_older
- gdp_per_capita
- extreme_poverty
- cardiovasc_death_rate
- diabetes_prevalence
- hospital_beds_per_thousand
- life_expectancy
- human_development_index

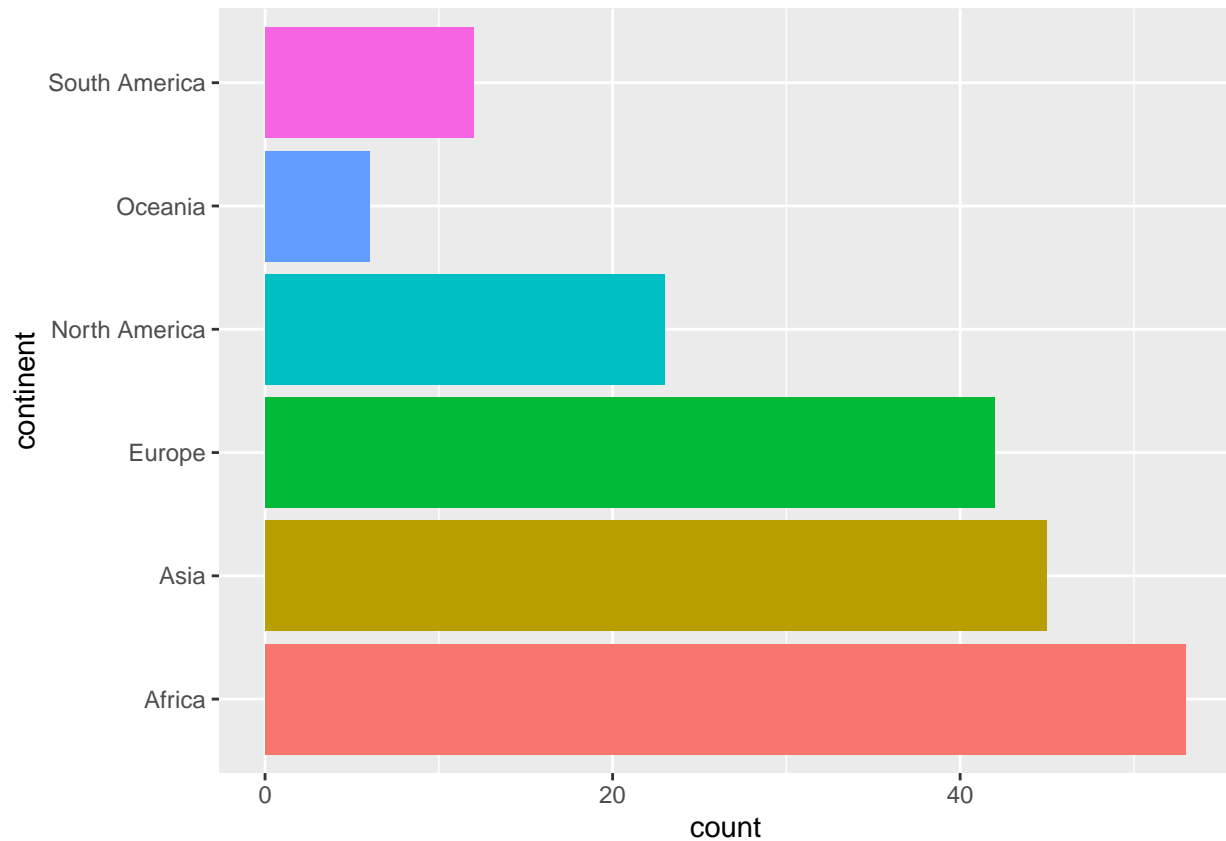
We select variables that we consider interesting to visualize, as the ones we haven't selected might be related to these or even ratios of them (in the case of total cases per million)

```
categorical <- c('location', 'continent', 'development')
interesting_vars <- c('total_cases', 'new_cases', 'total_deaths', 'stringency_index', 'population', 'populat
```

Plots with categorical variables

Countries per continent in the dataset

```
ggplot(data=data) +  
  geom_bar(aes(fill=continent, y=continent), show.legend = FALSE)
```

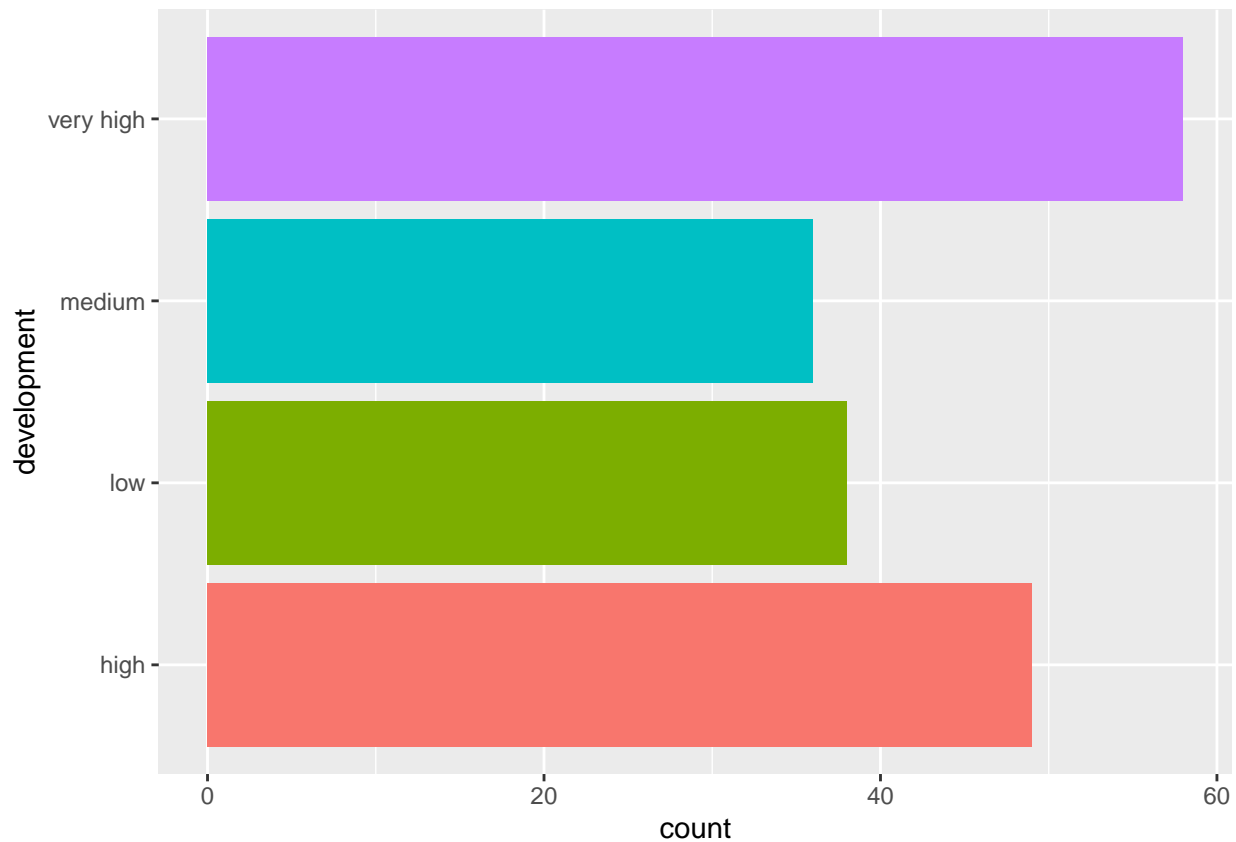


As a heads up, we will have more data points for continents with more countries, as there is an entry per country in the dataset. Therefore Africa will be our continent with the most entries, followed by Asia and Europe.

North america contains all of Central America and the Caribbean. Oceania contains some island countries and archipelagos in the pacific that belong to the continent but it is by far our continent with the least entries.

Amount of countries per HDI

```
ggplot(data=data) +  
  geom_bar(aes(fill=development, y=development), show.legend = FALSE)
```



For development (which is a variable constructed from the *human_development_index* variable), we have an even amount of countries per HDI, with very high development countries being the largest group.

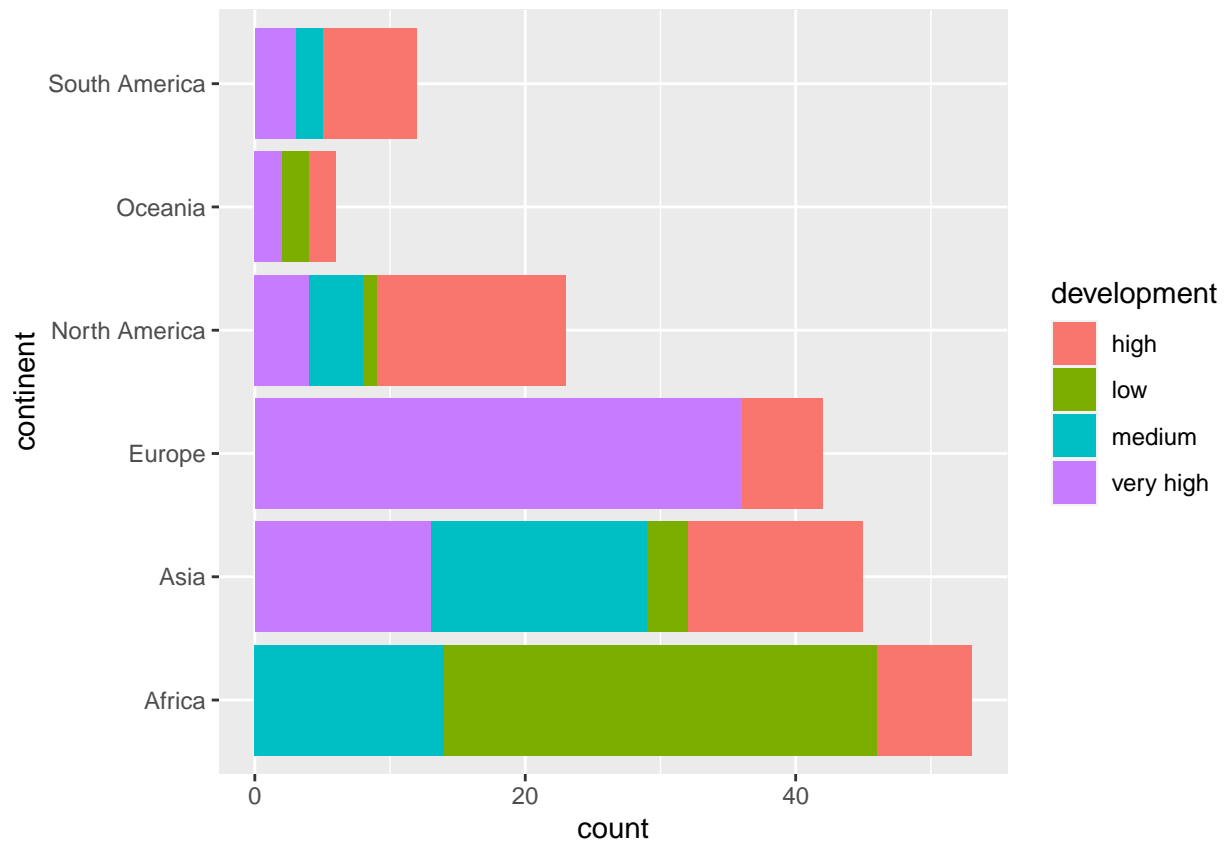
The *development* variable was constructed as follows from the *human_development_index* variable:

- *very high* for HDI of 0.800 and above
- *high* from 0.700 to 0.799
- *medium* from 0.550 to 0.699
- and *low* below 0.550.

We followed Wikipedia's criteria for the construction of this variable as it accurately represents and summarizes well HDI in 4 categories.

Countries per continent per HDI

```
ggplot(data=data) +  
  geom_bar(aes(fill=development ,y=continent))
```



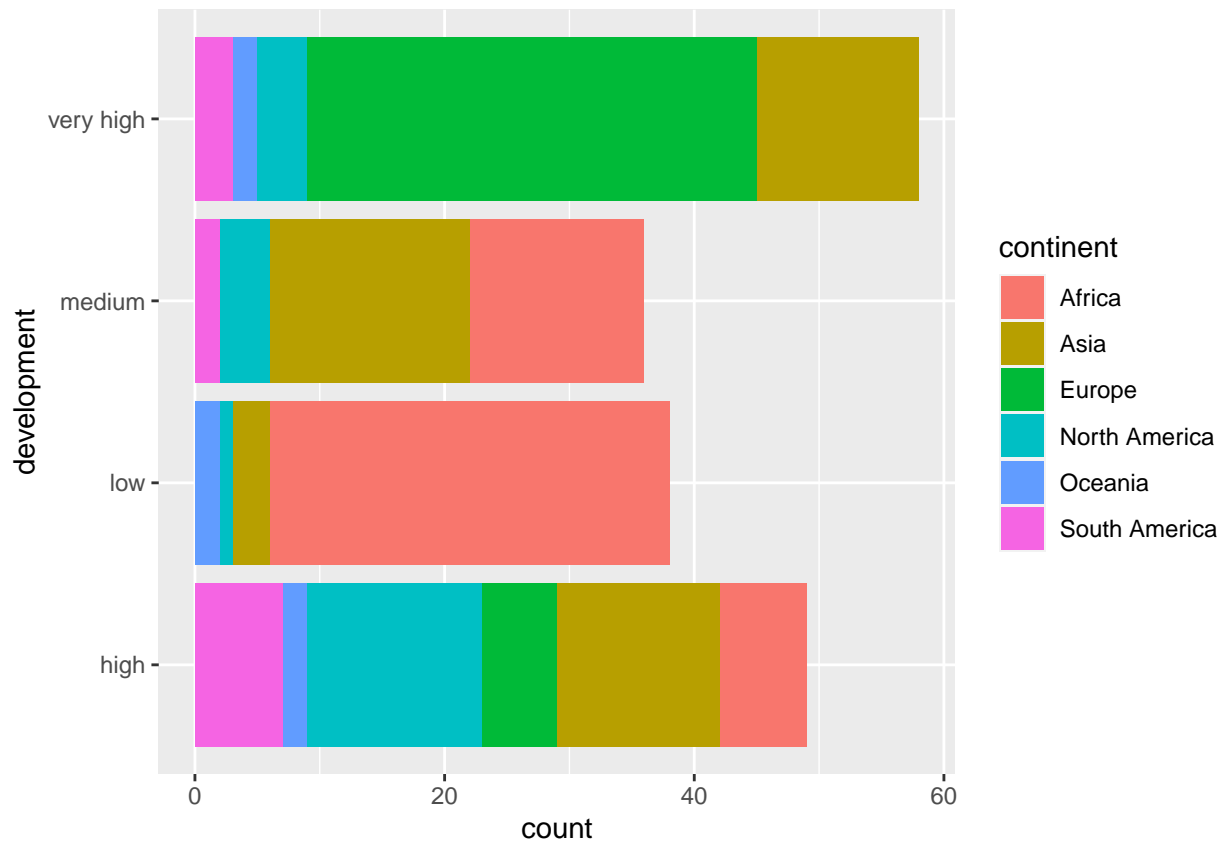
We can see here how many countries per continent correspond to which HDI using our constructed development variable. We can see Europe is fully composed of high and very high development countries. North America has most of its countries being high development countries, with a significant group of medium to low development countries, possibly located in Central America and the Caribbean. The very high development countries correspond to Bahamas, Barbados, Canada and the United States.

Africa is mostly composed of medium to low development countries and Asia is a bit of a tossup between high, very high and medium development countries.

To be fair, a country having high HDI does not imply that it isn't a developing country, as most countries with high HDI are, indeed, developing countries. This is simply a categorical representation of Human Development Index.

Proportions of HDI per continent

```
ggplot(data=data) +  
  geom_bar(aes(fill=continent, y=development))
```



Looking at the opposite plot to the previous one we can see that most of the very high development countries are located in Europe and Asia. With not a single very high development country in Africa and with most of the low development countries located in Africa.

Medium development countries are mostly located in Asia, Africa and the Americas with not a single medium or low development country located in Europe.

Which again, does not suggest at all that perhaps poverty situations do not exist in the European continent. For instance, about 30% of the population of Albania live under 5.50 USD a day. This figure corresponds to nearly a third of the population and this particular trend shows up for other countries in the Balkan Peninsula along with Moldova and Romania (for more examples).

Plots with numerical variables

Function to plot quantitative variables

```
plots <- function(dataset ,col, type, density=TRUE, bins='default', xtick_angles='default') {
  var <- dataset %>% dplyr::select(col)
  if (bins == 'default') {bins = rep(10,3)}
  if (xtick_angles == 'default') {xtick_angles = rep(90,3)}
  if (type == 'boxplot') {
    p1 <- dataset %>% ggplot(aes(x=var[,1])) +
      geom_boxplot() +
      ggtitle(str_interp("${col}")) +
      theme(axis.title.x=element_blank(),
            axis.text.y=element_blank())
    p2 <- dataset %>% ggplot(aes(x=var[,1], fill=continent)) +
      geom_boxplot() +
      ggtitle(str_interp("${col} grouped by continent")) +
      theme(axis.title.x=element_blank(),
            axis.text.y=element_blank())
    p3 <- dataset %>% ggplot(aes(x=var[,1], fill=development)) +
      geom_boxplot() +
      ggtitle(str_interp("${col} grouped by development")) +
      theme(axis.title.x=element_blank(),
            axis.text.y=element_blank())
  } else if (type == 'hist') {
    p1 <- dataset %>% ggplot(aes(x=var[,1])) +
      geom_histogram(aes(y=..density..), bins=bins[1]) +
      geom_density() +
      ggtitle(str_interp("${col}")) +
      theme(axis.title.x=element_blank(),
            axis.text.x = element_text(angle = xtick_angles[1]))
    if (density == FALSE) {
      p2 <- dataset %>% ggplot(aes(x=var[,1], fill=continent)) +
        geom_histogram(show.legend = FALSE,bins=bins[2]) +
        ggtitle(str_interp("${col} by continent")) +
        theme(axis.title.x=element_blank(),
              axis.text.x = element_text(angle = xtick_angles[2])) +
        facet_wrap(~continent, nrow = 1)
      p3 <- dataset %>% ggplot(aes(x=var[,1], fill=development)) +
        geom_histogram(show.legend = FALSE,bins=bins[3]) +
        ggtitle(str_interp("${col} by development")) +
        theme(axis.title.x=element_blank(),
              axis.text.x = element_text(angle = xtick_angles[3])) +
        facet_wrap(~development, nrow = 1)
    } else {
      p2 <- dataset %>% ggplot(aes(x=var[,1], fill=continent)) +
        geom_histogram(show.legend = FALSE,bins=bins[2],aes(y=..density..)) +
        geom_density(show.legend = FALSE) +
        ggtitle(str_interp("${col} by continent")) +
        theme(axis.title.x=element_blank(),
              axis.text.x = element_text(angle = xtick_angles[2])) +
        facet_wrap(~continent, nrow = 1)
      p3 <- dataset %>% ggplot(aes(x=var[,1], fill=development)) +
        geom_histogram(show.legend = FALSE,bins=bins[3],aes(y=..density..)) +
```



```

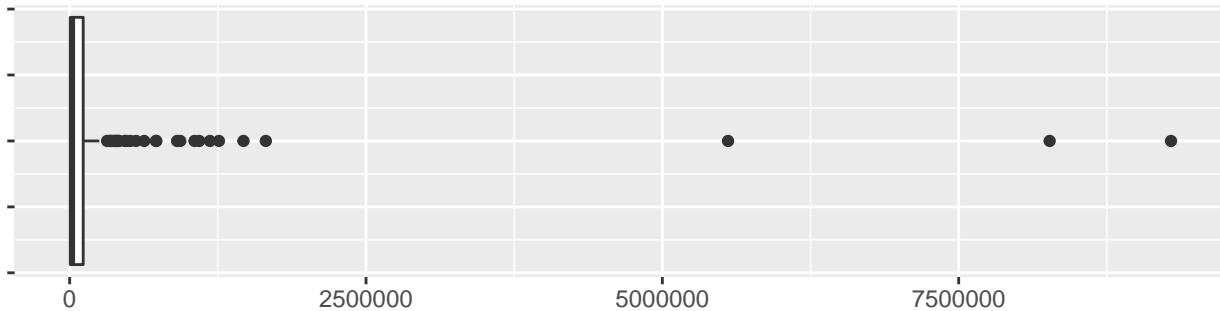
    geom_density(show.legend = FALSE) +
    ggtitle(str_interp("${col} by development")) +
    theme(axis.title.x=element_blank(),
           axis.text.x = element_text(angle = xtick_angles[3])) +
    facet_wrap(~development, nrow = 1)
  }
}
grid.arrange(p1,p2,p3, nrow=3)
}

```

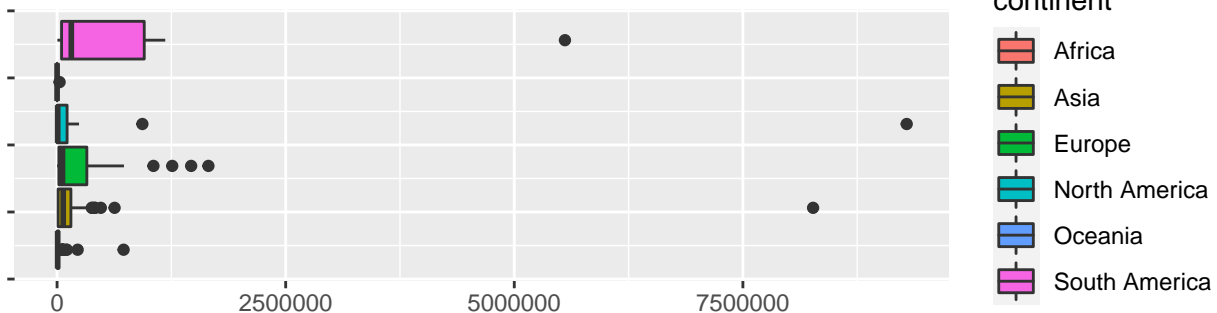
Boxplots for total cases of COVID-19

```
plots(dataset=data, col='total_cases',type='boxplot')
```

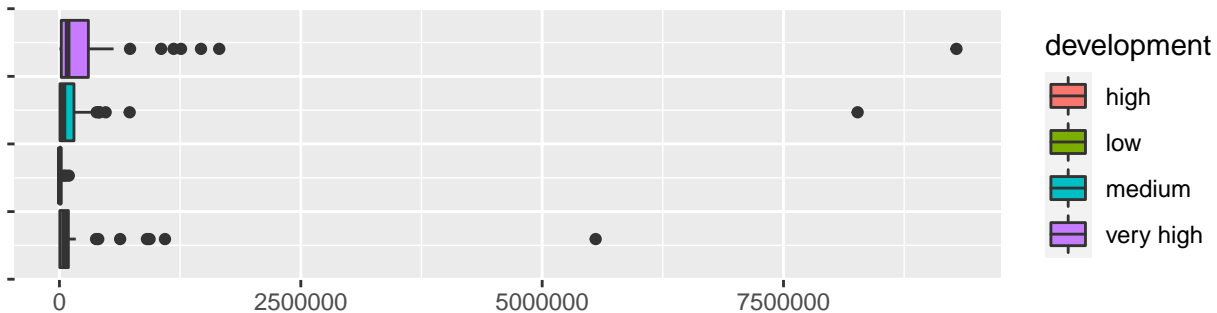
total_cases



total_cases grouped by continent



total_cases grouped by development

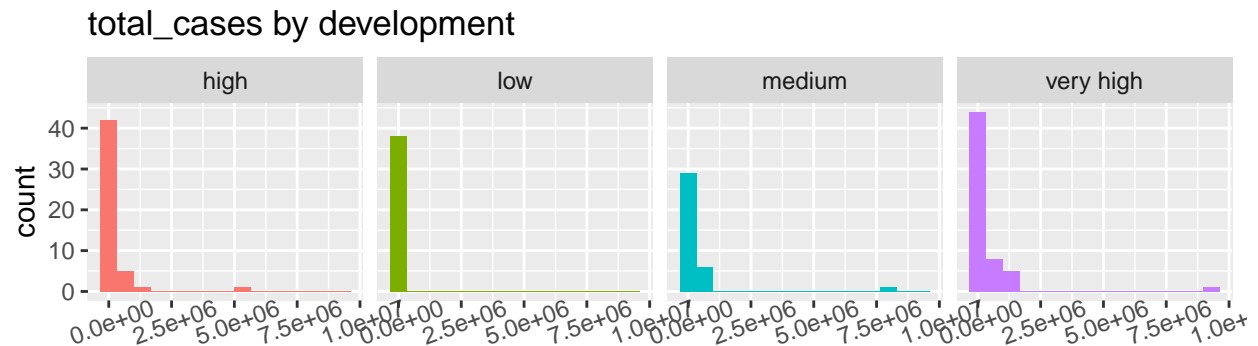
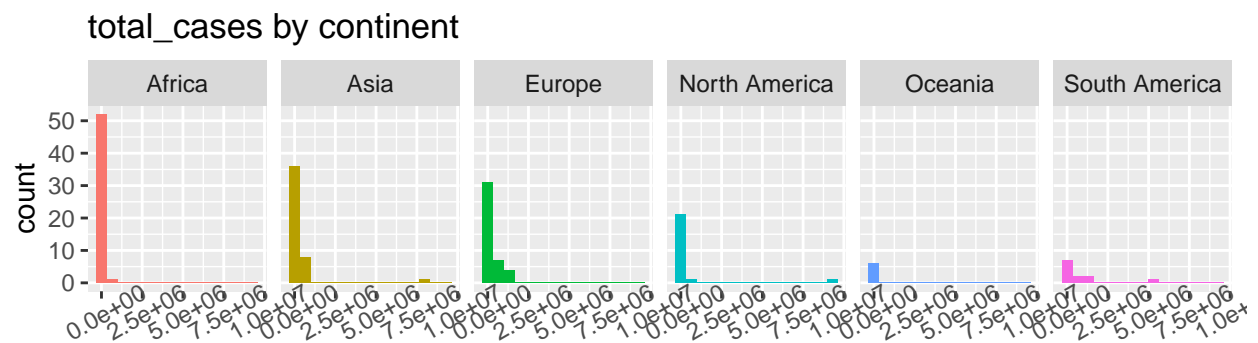
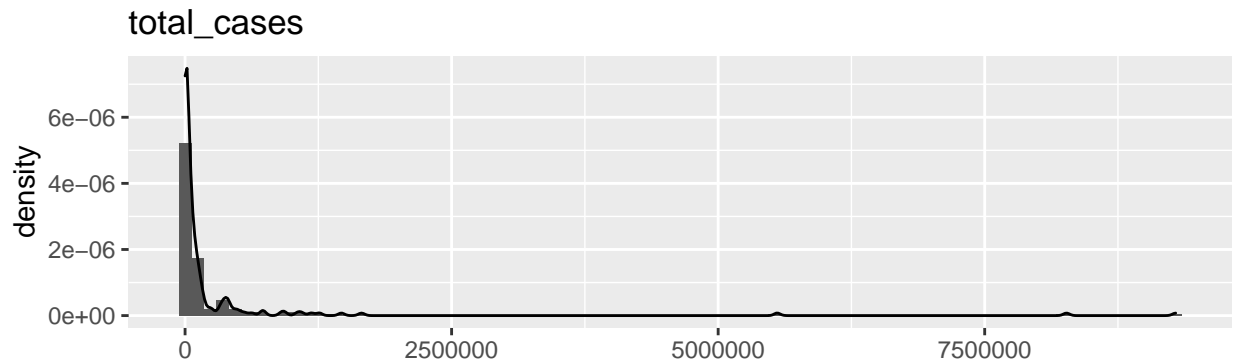


On the first box-plot of the variable, we can observe that the distribution of it is very right skewed, with some outliers. And the box-plots grouped by continents tell us that the country that has the most of the total cases is the US from continent North America which has a very high HDI. Meanwhile, the country that has the most of the total cases of Asia is India (second of the world), and has medium HDI. The third country that has the most of the cases is Brazil from South America with high HDI.

We can probably say that these three countries are the outliers for the variable total cases.

Histogram and kernel density for total cases of COVID-19

```
plots(dataset=data, col='total_cases',type='hist', density=FALSE, bins = c(80,15,15),xtick_angles=c(0,30,45))
```

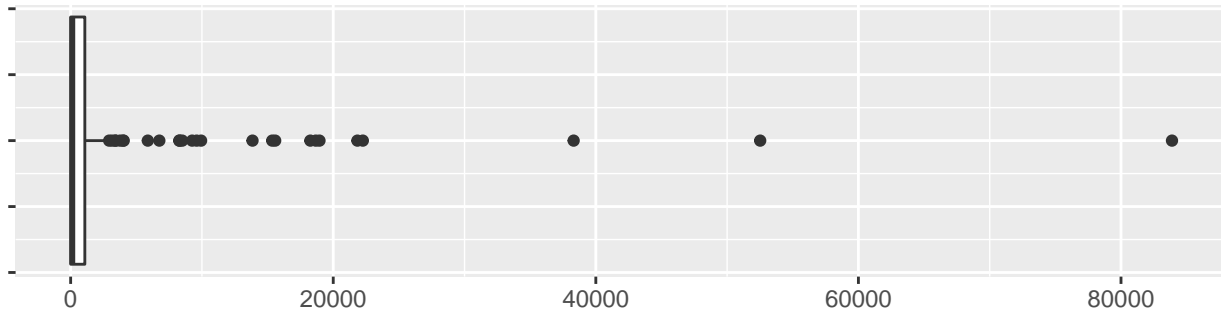


Observing this graph, we can confirm that the distribution is very right-skewed. Only Africa, Europe and Oceania don't have outliers. But it is probably because we don't have the dataset updated yet (we have the data-set updated on 3rd of November, 2020). About the development of different countries, we can't group the countries in terms of how they have developed by the total cases of COVID-19 they have.

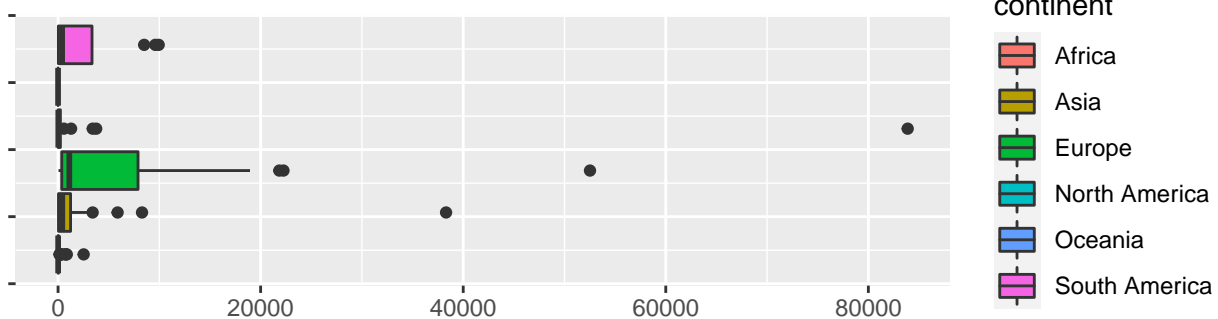
Boxplots for new cases of COVID-19

```
plots(dataset=data, col='new_cases',type='boxplot')
```

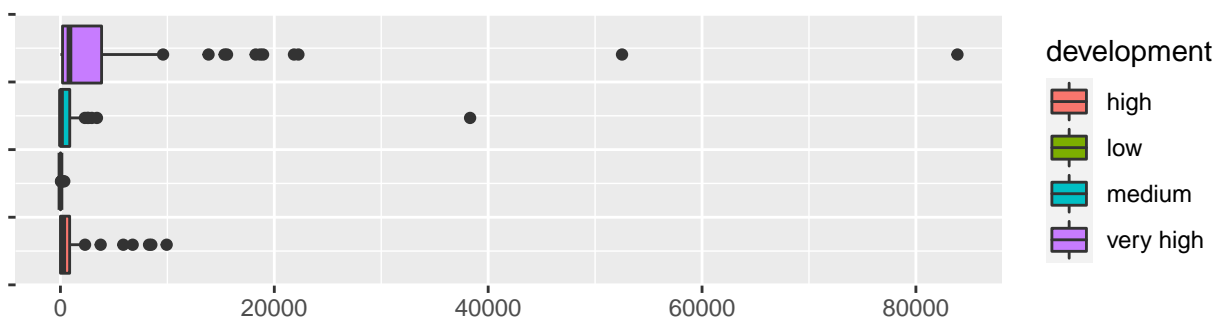
new_cases



new_cases grouped by continent



new_cases grouped by development



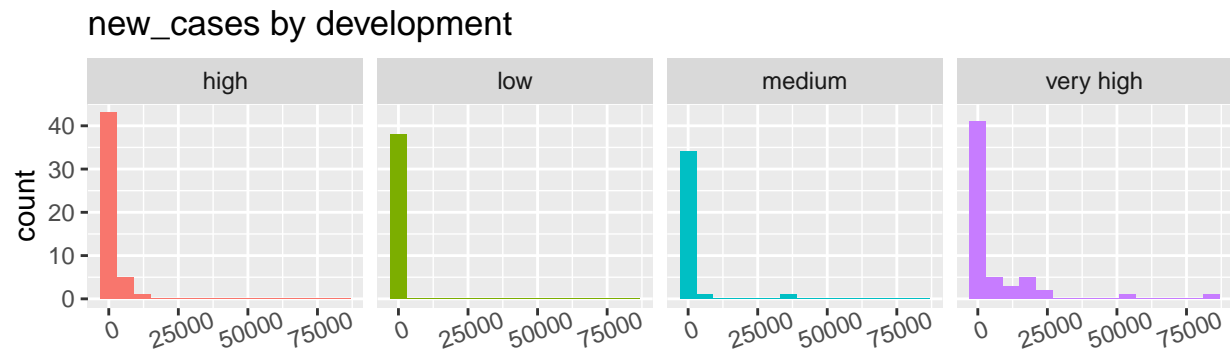
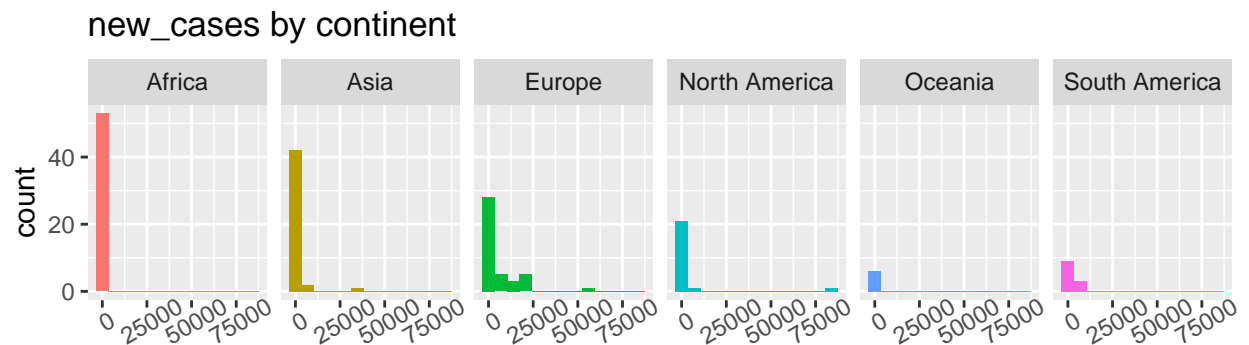
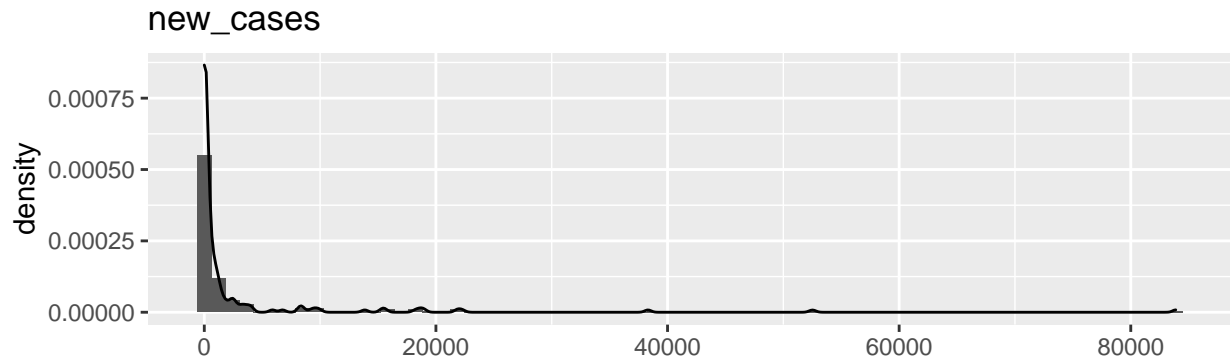
In the first box-plot above, we see that the distribution is very right skewed with some outliers. The country that has the most of the new cases is the US.

Observing the box-plot of new cases grouped by continent it is obvious that the country of North America that has the most of the new cases is the US. And the second country that has the most of the new cases is in Europe, France. Both of them have a very high Human Development Index. The third country that has the most of the new cases is India from Asia with medium HDI.

Another thing to mention is that the countries that have the most of the new cases is very related with the previous variable, which is the total cases, they have similar characteristics.

Histogram and kernel density for new cases of COVID-19

```
plots(dataset=data, col='new_cases',type='hist', density=FALSE, bins = c(70,13,15),xtick_angles=c(0,30,45))
```



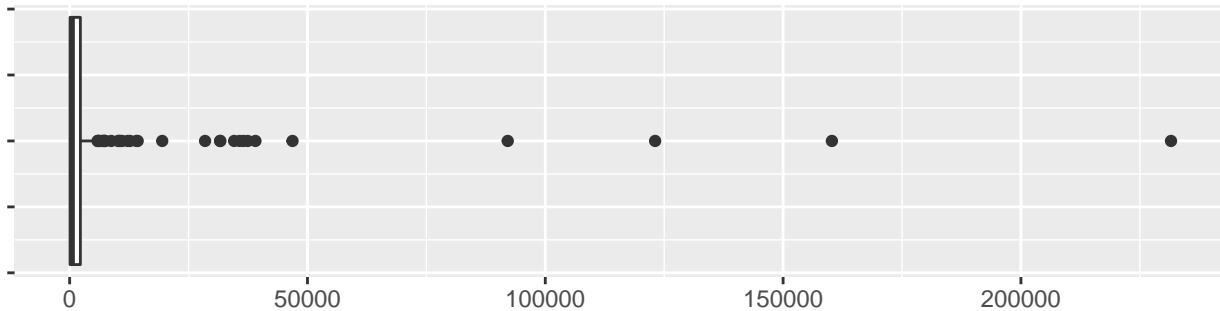
Observing this graph, we can confirm that the distribution is very right-skewed with some outliers (especially the US, France and India).

In the histograms of new cases by continent, we can see that Europe has the most dispersed distribution comparing to other continents, which means that the countries of Europe have very different values of new cases from each other. About the development of different countries, we can't group the countries in terms of how they have developed by the new cases per day of COVID-19 they have, the distribution of countries that have a very high Human development Index have a some outliers.

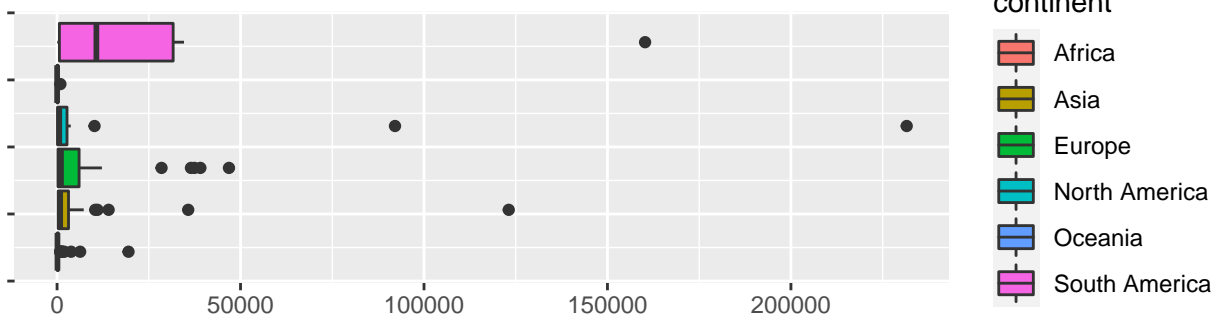
Boxplots for total deaths due to COVID-19

```
plots(dataset=data, col='total_deaths',type='boxplot')
```

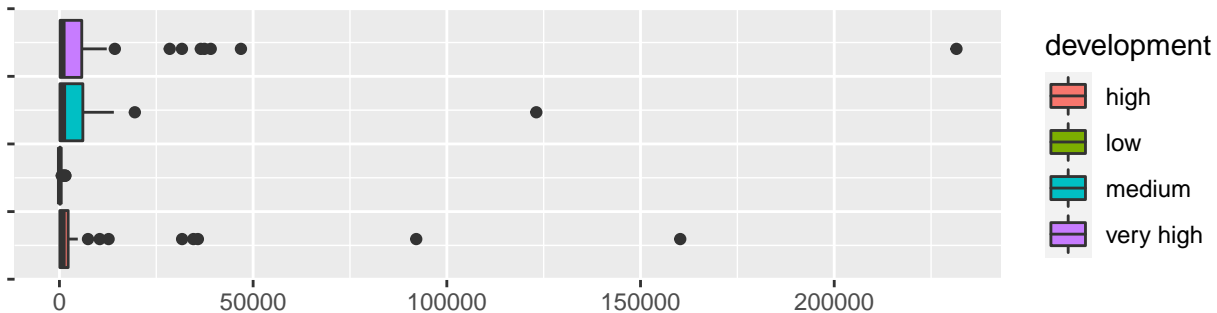
total_deaths



total_deaths grouped by continent



total_deaths grouped by development

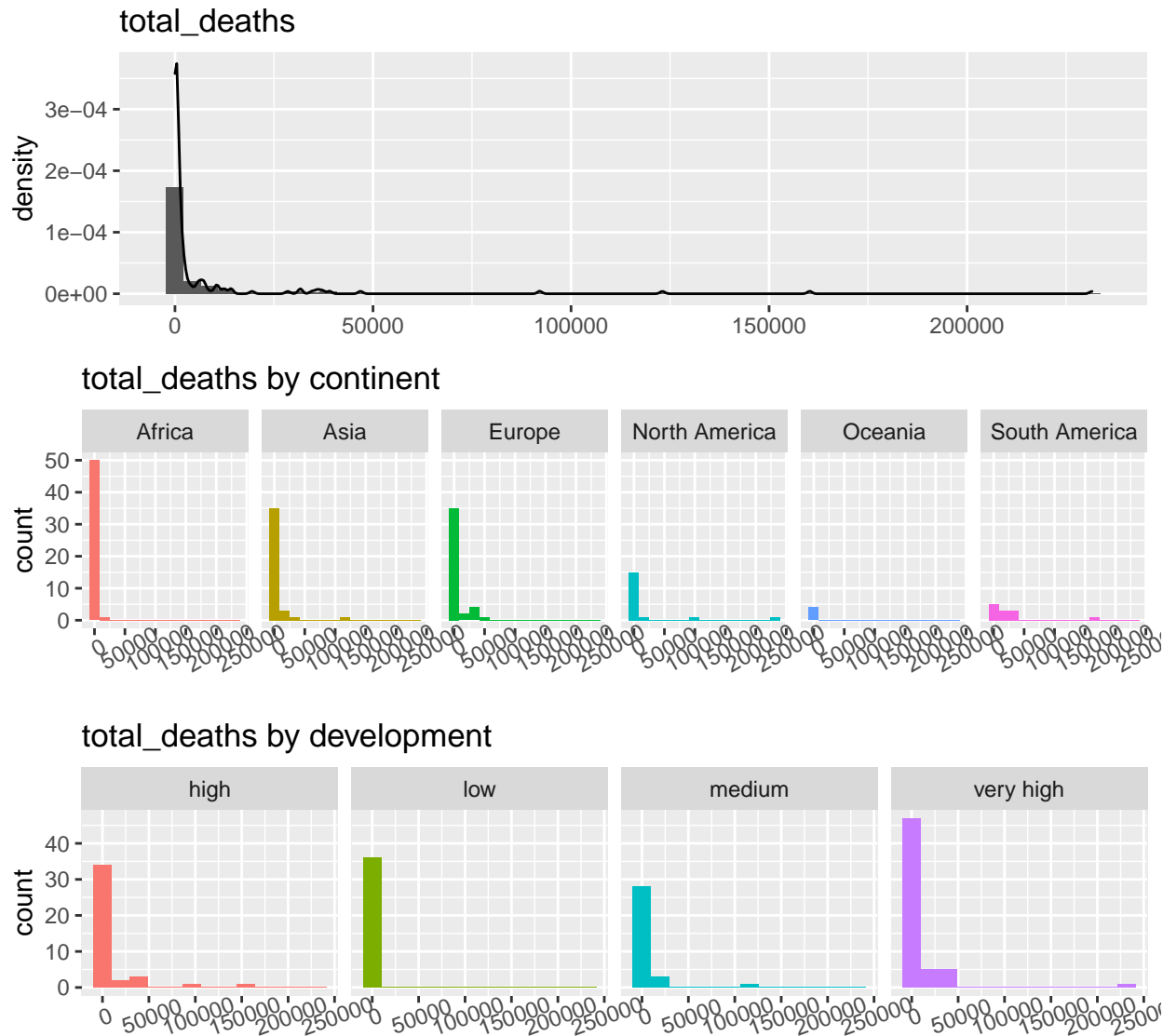


From the box-plots above we can say that this variable of new deaths is very likely distributed with the variables total cases and new cases. These three variables are all right-skewed and all have some outliers.

In this case, the country that has the most of the total deaths is still the US. It is obvious that the country of North America that has the most of the total deaths is the US. And the second country that has the most of the total deaths is in South America, Brazil. The third country that has the most of the total deaths is from Asia, India.

Histogram and kernel density for total deaths due to COVID-19

```
plots(dataset=data, col='total_deaths',type='hist', density=FALSE, bins = c(55,15,13),xtick_angles=c(0,1
```



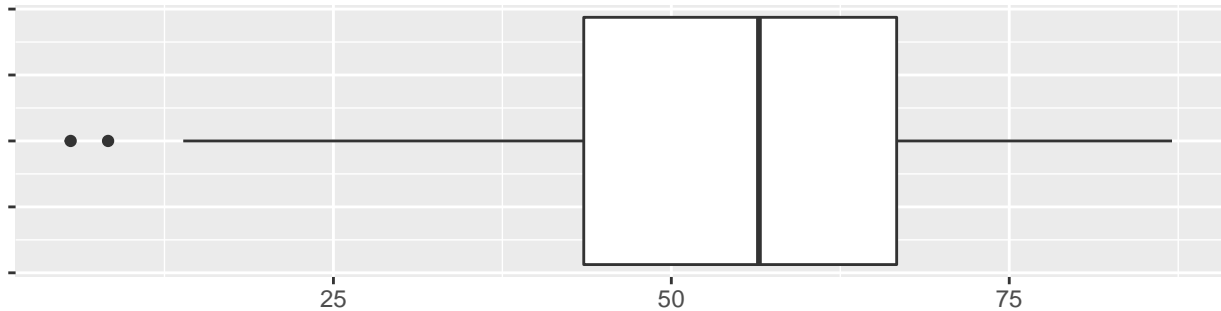
Again, the distribution is very right-skewed with some outliers (the US, Brazil and India).

In the histograms of new cases by continent and histograms by development, the distributions are also very right-skewed, some of them have outliers. We still can't group the countries in terms of how they have developed by the total deaths of COVID-19 they have, the distribution of countries that have a very high, high and medium Human development Index have a some outliers.

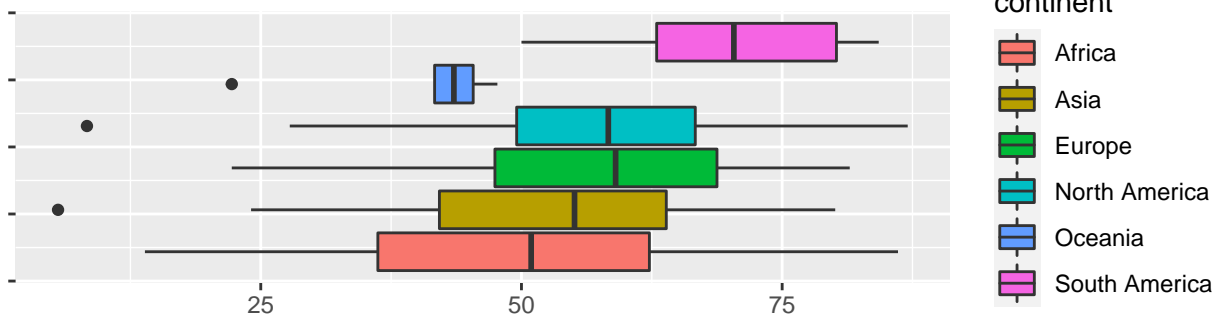
Boxplots for stringency index (how strict measures are)

```
plots(dataset=data, col='stringency_index', type='boxplot')
```

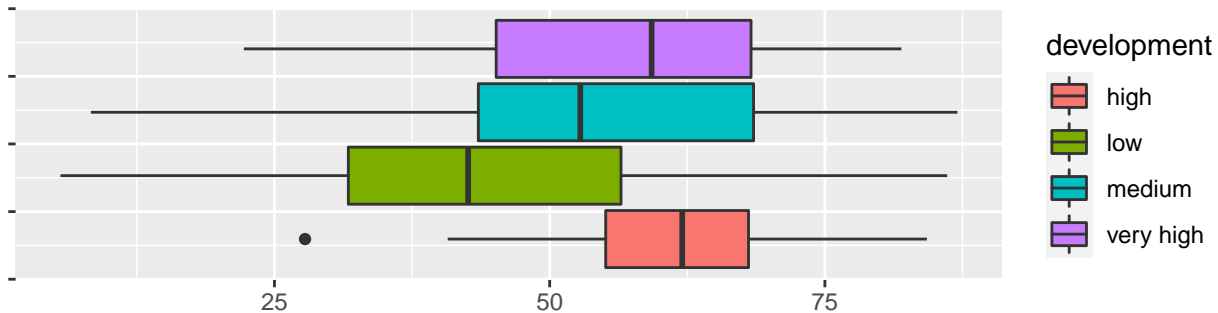
stringency_index



stringency_index grouped by continent



stringency_index grouped by development

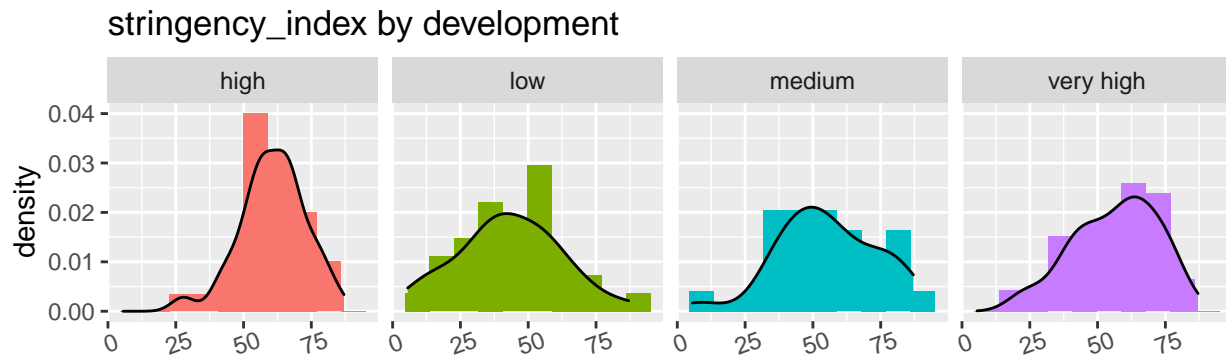
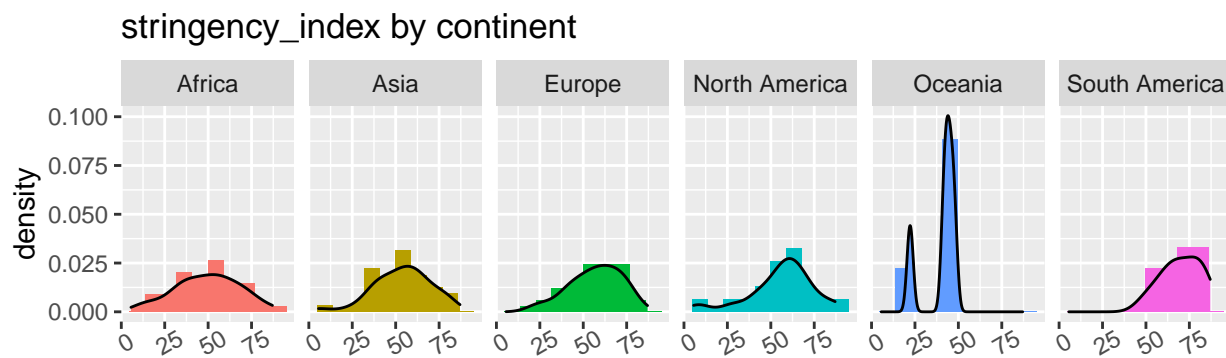
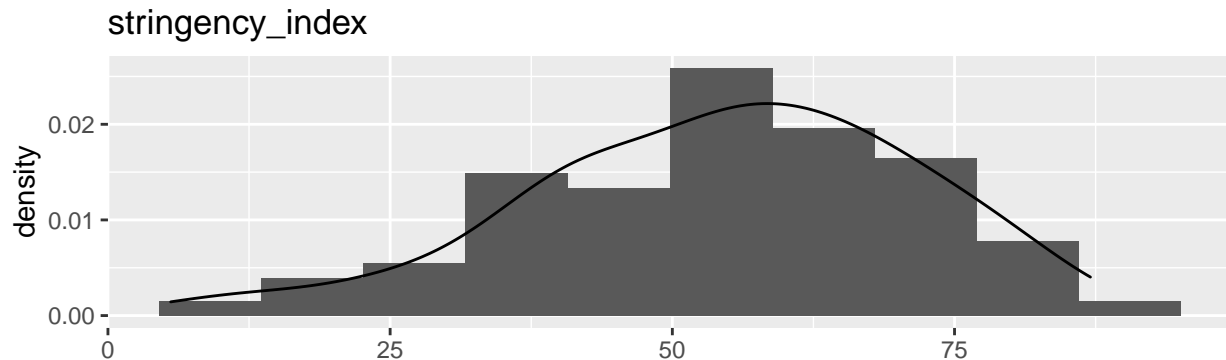


From the box-plot above we can say that the global distribution of stringency index is a little left skewed, but it is the most normally distributed until now. There are some countries that have really low stringency index, for example Afghanistan from Asia has stringency index equals to 5.56, or Nicaragua from North America with stringency index 8.33.

Observing the box-plot of stringency index grouped by continent we can see that South America has the most strict measurements and Oceania has the least strict measurements. The rest of them have similar distribution on stringency index. Look at the stringency index grouped by development we notice that the countries which have low HDI have less stringency index (using the criterion of quantiles).

Histogram and kernel density for stringency index

```
plots(dataset=data, col='stringency_index',type='hist', density=TRUE,xtick_angles=c(0,30,20))
```

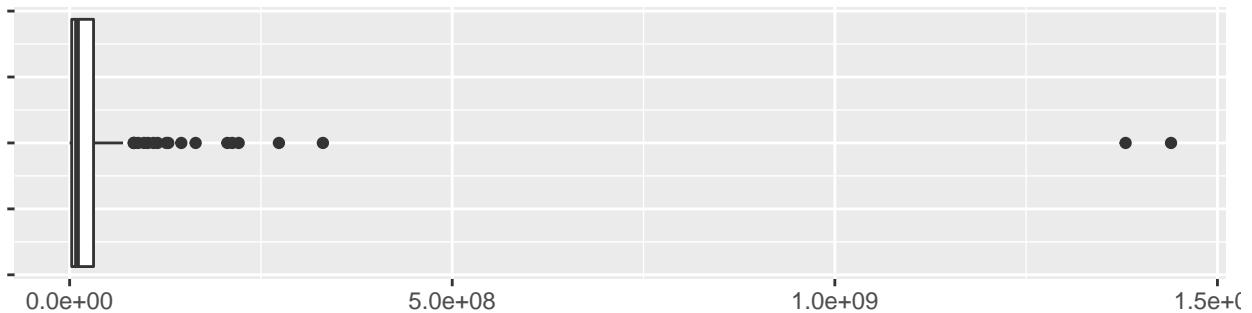


The distribution is quite symmetric distributed for different continents. Except for Oceania, there are some countries have really low stringency index; and South America, the distribution of this variable is quite left skewed. We can probably distinguish the countries with low HDI from others, these countries usually have less stringency index.

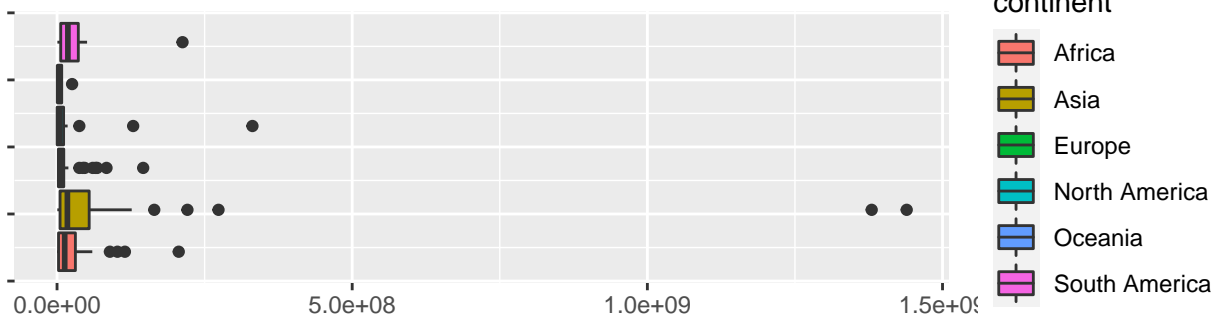
Boxplots for population

```
plots(dataset=data, col='population',type='boxplot')
```

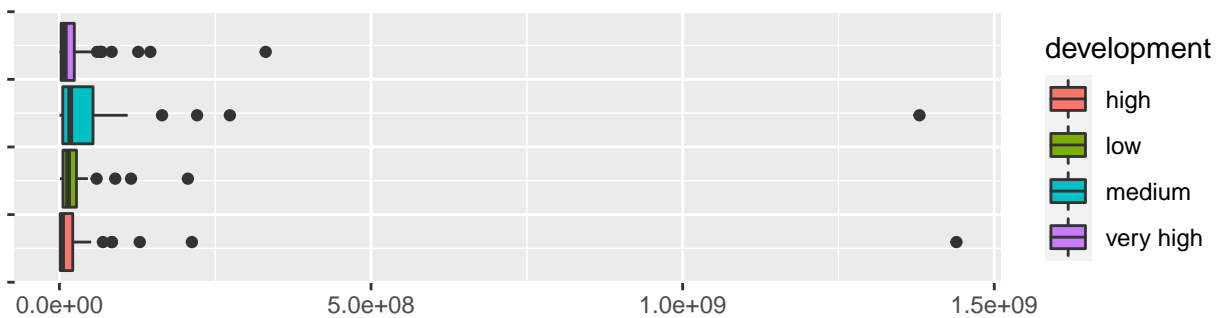
population



population grouped by continent



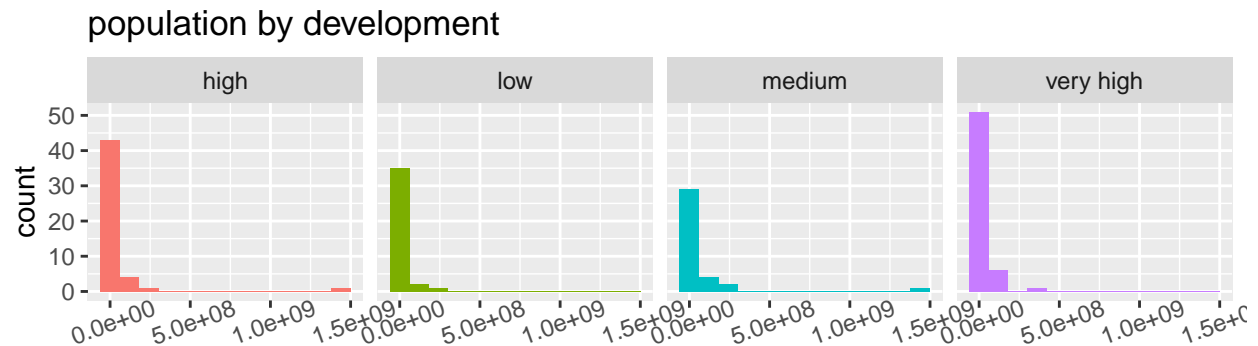
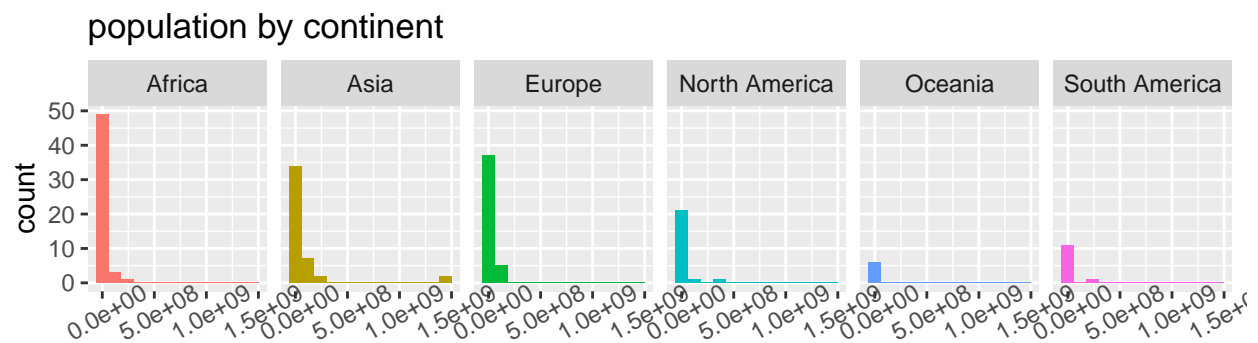
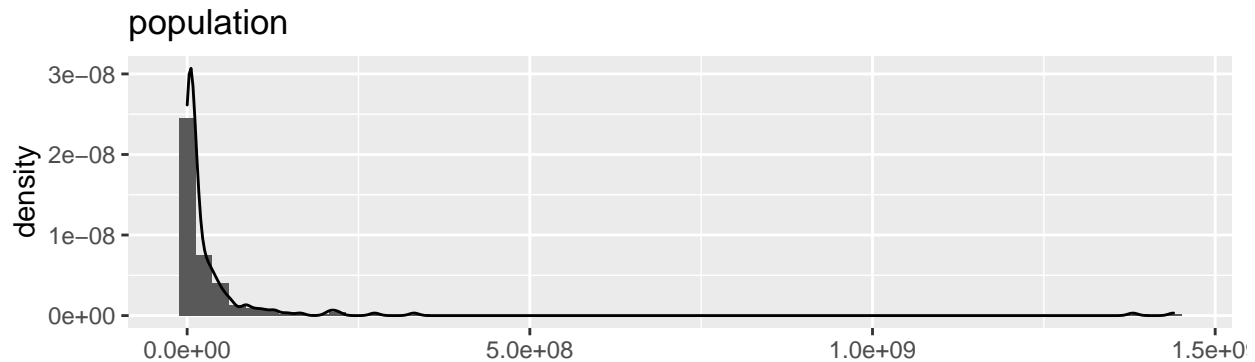
population grouped by development



Observing the box-plot above we can see that the distribution of population is very right skewed, some countries have way more population than others. For example, China has the most population of all, India is the second, these two countries are most exaggerated outliers from the plot.

Histogram and kernel density for population

```
plots(dataset=data, col='population',type='hist', density=FALSE, bins = c(60,13,13),xtick_angles=c(0,30
```

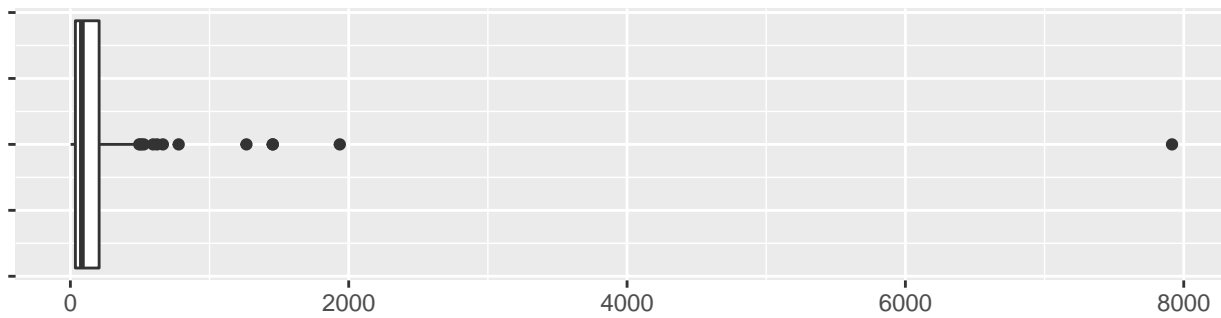


The distribution is very right-skewed, the population of each country is very different from others, but the variable population does not provide any information of whether the country has high HDI or not.

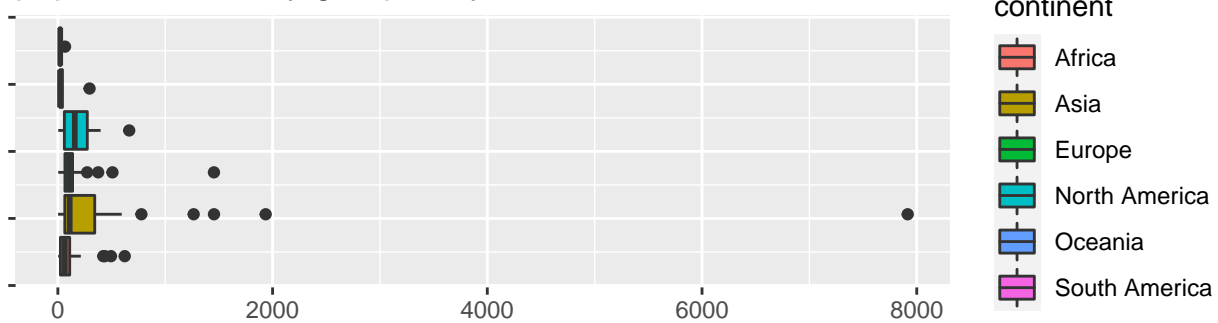
Boxplots for population density

```
plots(dataset=data, col='population_density',type='boxplot')
```

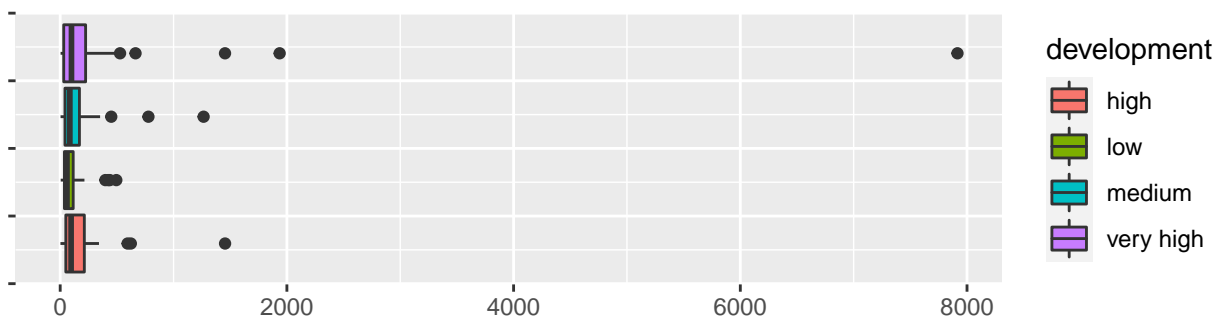
population_density



population_density grouped by continent



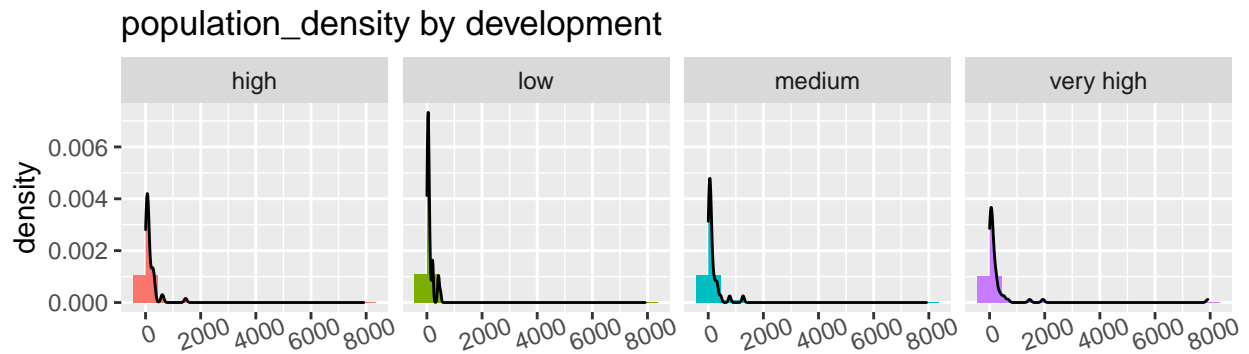
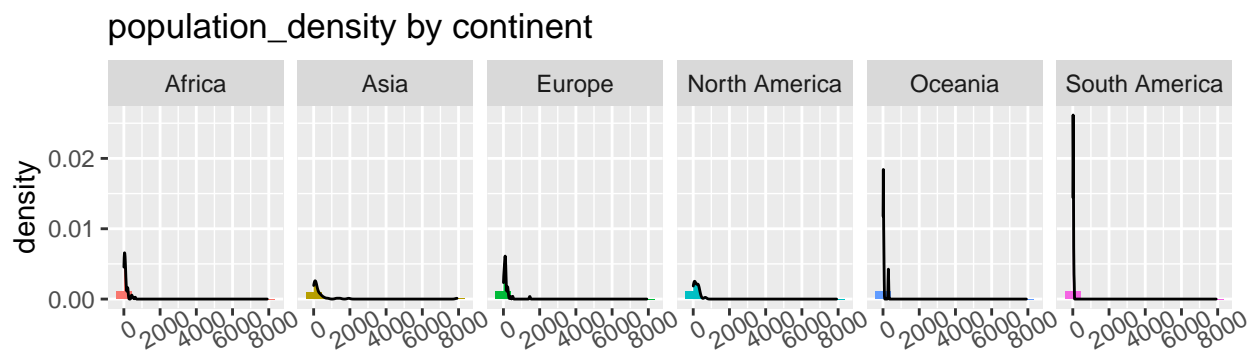
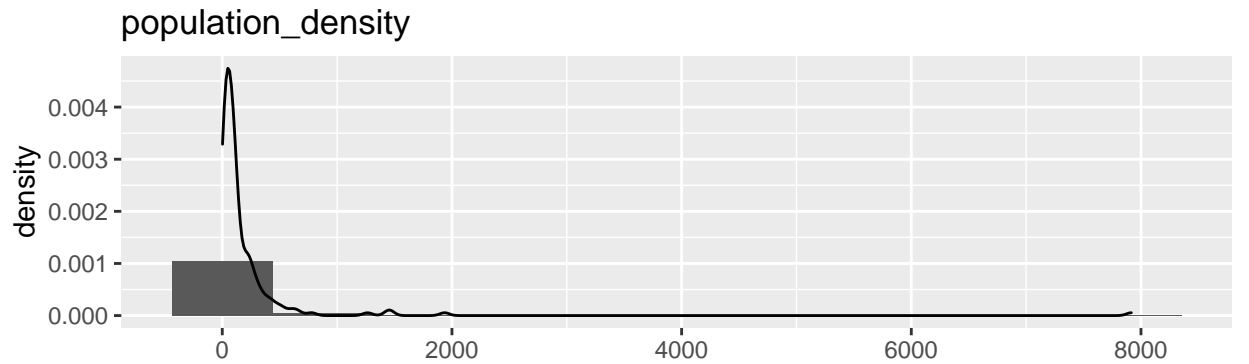
population_density grouped by development



Population density is a measurement of population per unit area. Observing the previous box-plot above we can see that the distribution of population density is very likely distributed as population distribution, it is very right skewed, some countries have really high population density. For instance, Singapore has the most population density of all with a value of 7915.731, it is a small country of Asia with very high HDI.

Histogram and kernel density for population density

```
plots(dataset=data, col='population_density',type='hist', density=TRUE, xtick_angles=c(0,30,20))
```

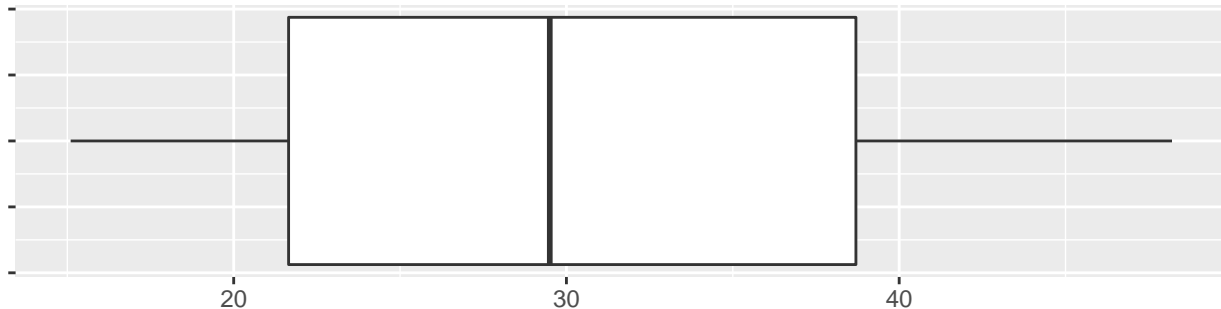


The distribution is very right-skewed, the population density of each country is very different from others. And the variable does not provide any information of whether the country has high HDI or not.

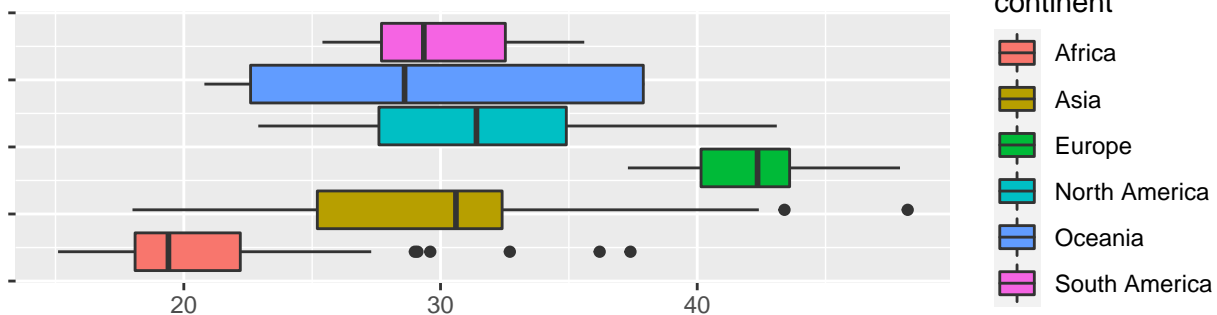
Boxplots for median age

```
plots(dataset=data, col='median_age',type='boxplot')
```

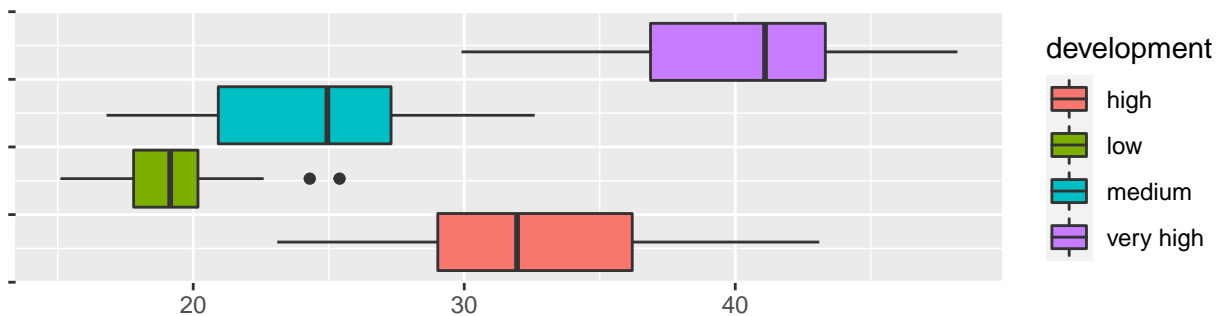
median_age



median_age grouped by continent



median_age grouped by development



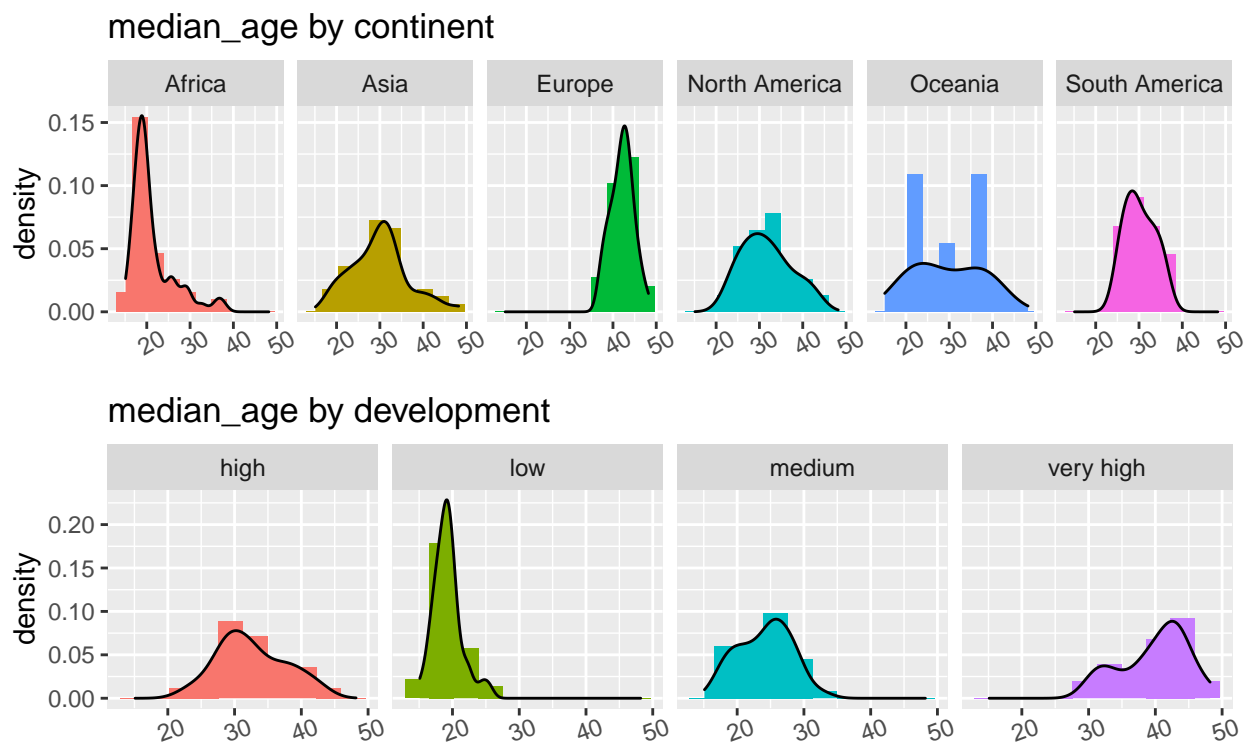
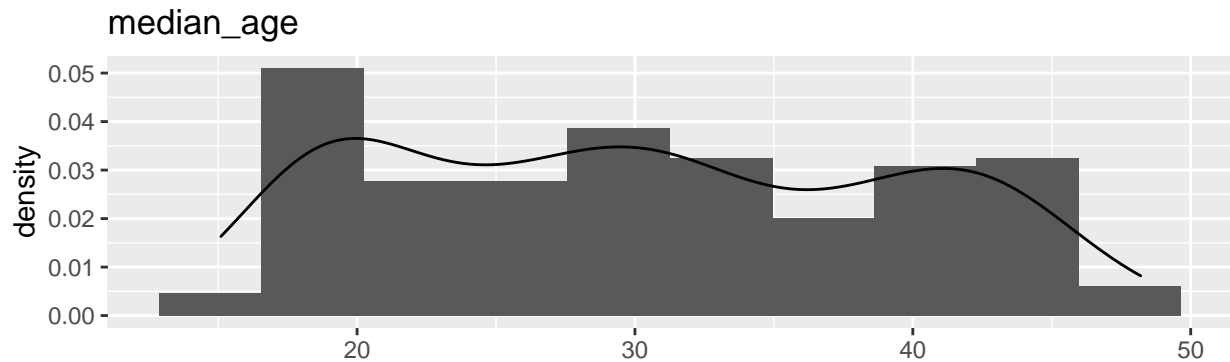
Observing the box-plot for the global median age we can notice that the distribution of it is quite symmetric. The majority of the median age of different countries is located between 20-40.

But from the grouped box-plots we can find something really interesting: - For the box-plots grouped by continent we can see that the median age of Europe is larger (more than 40) than the rest of the continents while Africa has the least median age (less than 20) with some “outliers” that have similar median age as other continents.

- For the box-plots grouped by development we detect that usually higher developed a country, larger the median age, e.g. the countries that have very high HDI have median of median age more than 40, and the countries that have low HDI have median of median age less than 20. From that, we can conclude that the majority of countries from Africa has low HDI while majority of countries from Europe has very high HDI.

Histogram and kernel density for median age

```
plots(dataset=data, col='median_age',type='hist', density=TRUE, xtick_angles=c(0,30,20))
```



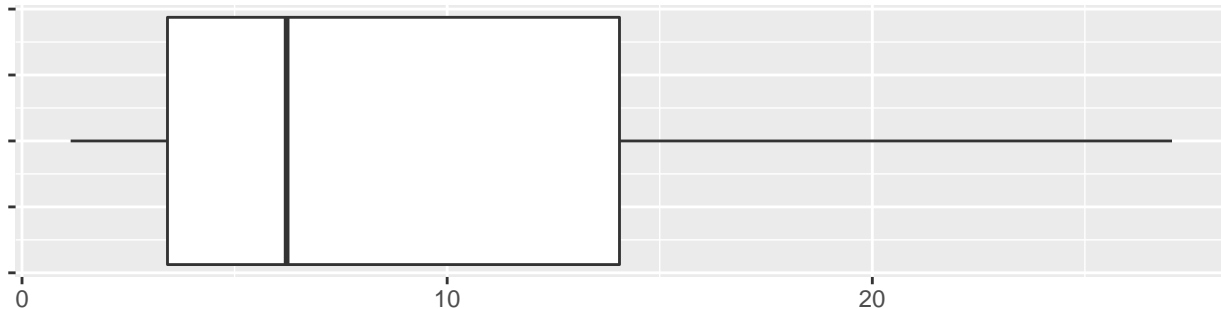
The distributions of median age for different continent are very different. The distribution of Africa is right skewed while others are symmetric.

The distributions of Asia, North America, Oceania and South America are more flat (platykurtic), and the distributions of Africa and Europe are more concentrated (leptokurtic).

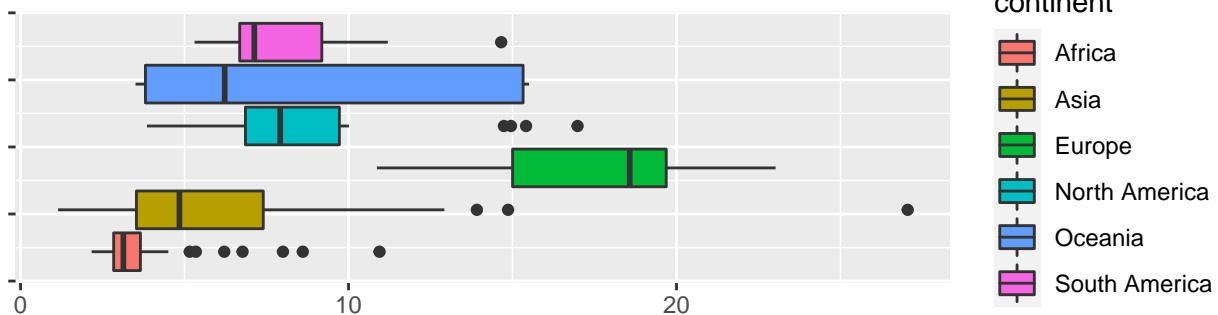
Boxplots for the percentage of population aged 65 or older

```
plots(dataset=data, col='aged_65_older', type='boxplot')
```

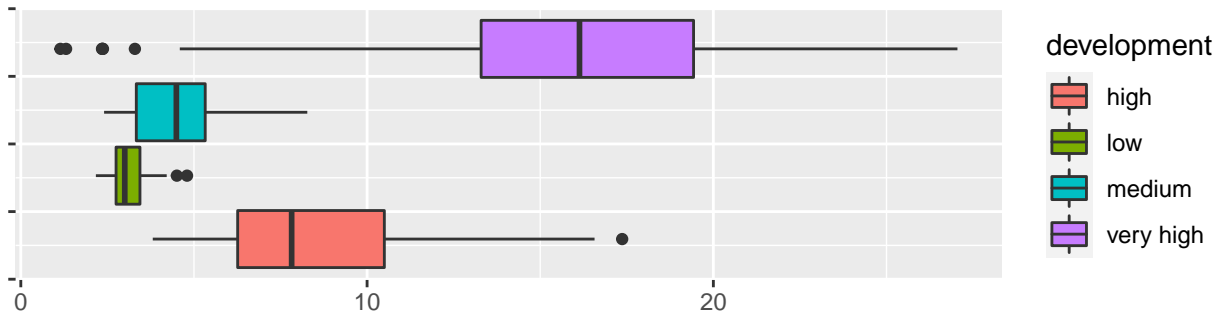
aged_65_older



aged_65_older grouped by continent



aged_65_older grouped by development



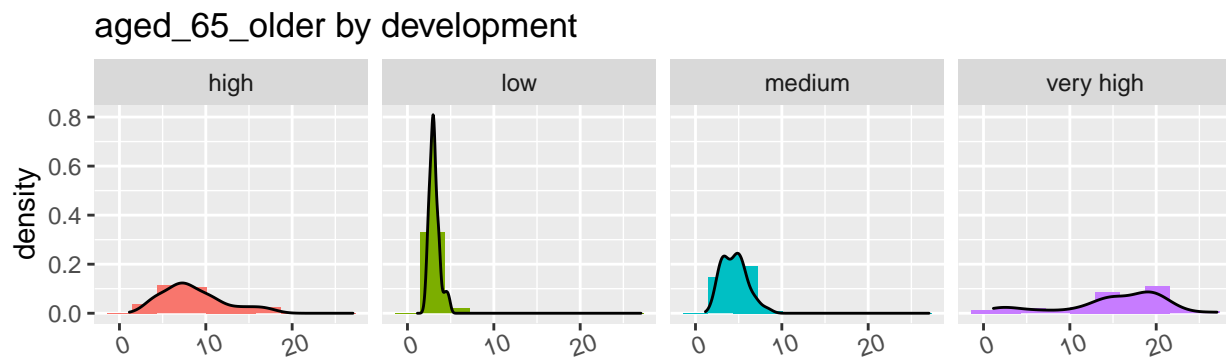
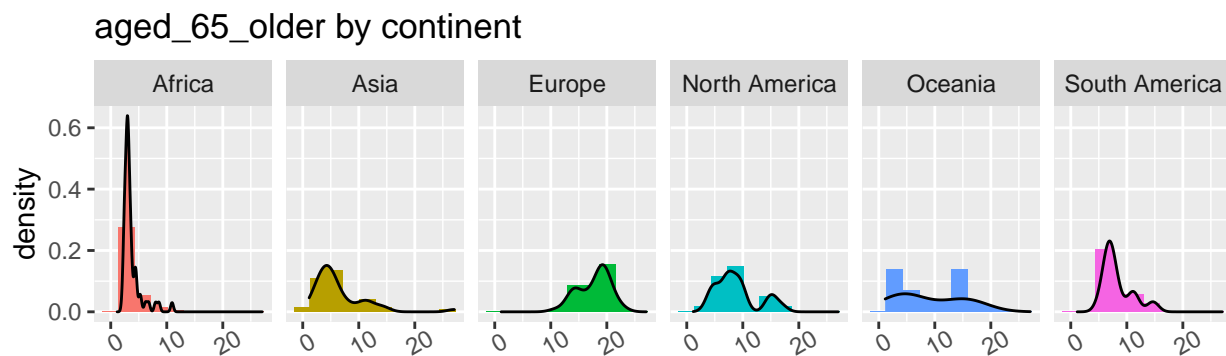
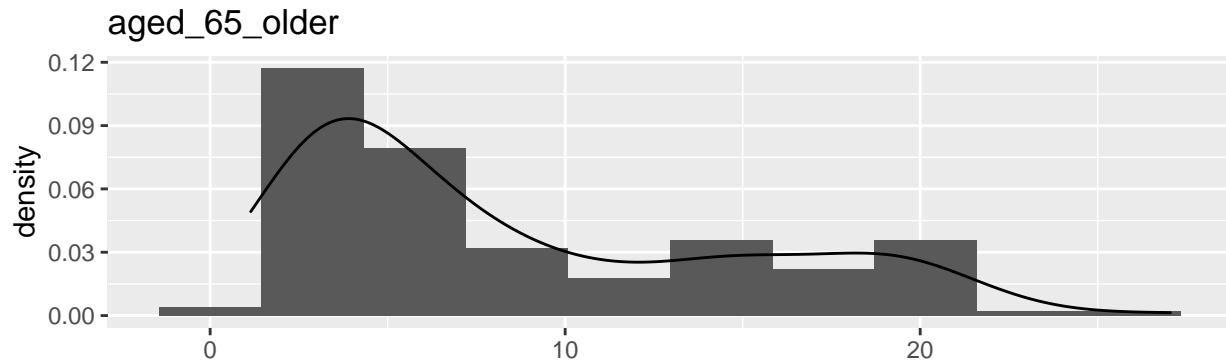
From the box-plot for the global percentage of population aged 65 or older we can notice that the distribution is right skewed. There are more than half of the countries have less than 10% of population of aged 65 or older.

However, the grouped box-plots give us more interesting results: - In the box-plots grouped by continent we can see that the percentage of population aged 65 or older of Europe is larger (more than 10%) than the rest of the continents while Africa has the least median age (generally less than 10%) with some “outliers” that have similar values as other continents.

- In the box-plots grouped by development we detect that usually higher developed a country, larger the median age. These box-plots kind of give us the same information as the median age of each country, although this variable is not as clear as the previous one, median age.

Histogram and kernel density for the percentage of population aged 65 or older

```
plots(dataset=data, col='aged_65_older',type='hist', density=TRUE, xtick_angles=c(0,30,20))
```



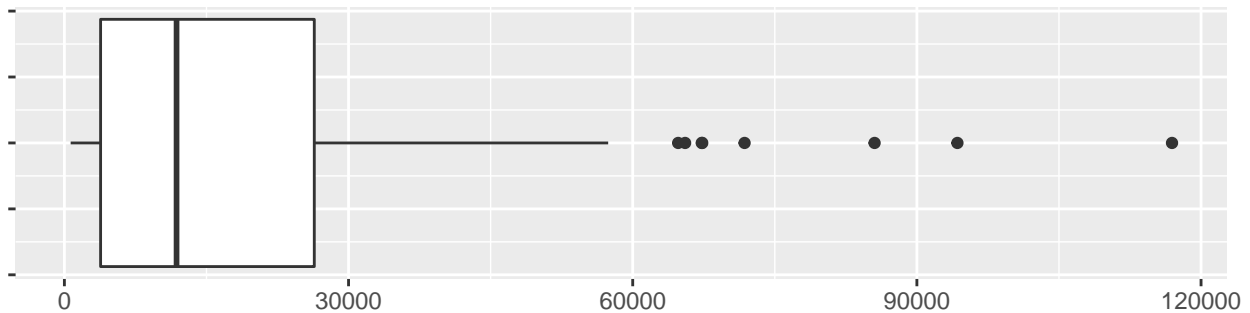
The distributions of this variable for different continent are very different. The distributions of Africa, Asia, North America and South America are right skewed while the distribution of Europe is left skewed.

The distribution of Africa is more concentrated (leptokurtic) while others are more flat. The distributions of countries that have low HDI are more concentrated, and they usually have less percentage of population of aged 65 or older.

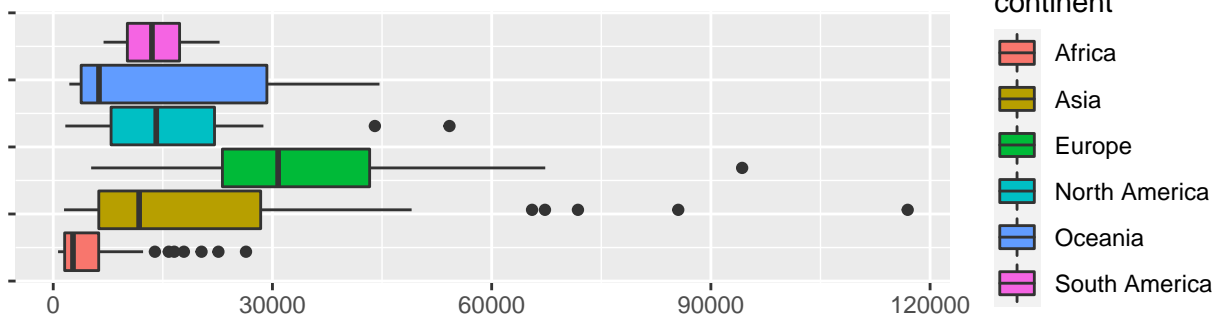
Boxplots for GDP per capita

```
plots(dataset=data, col='gdp_per_capita',type='boxplot')
```

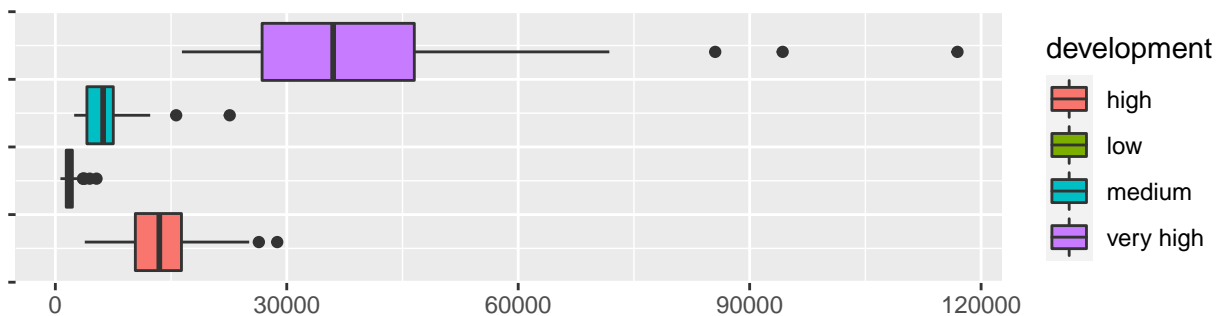
gdp_per_capita



gdp_per_capita grouped by continent



gdp_per_capita grouped by development



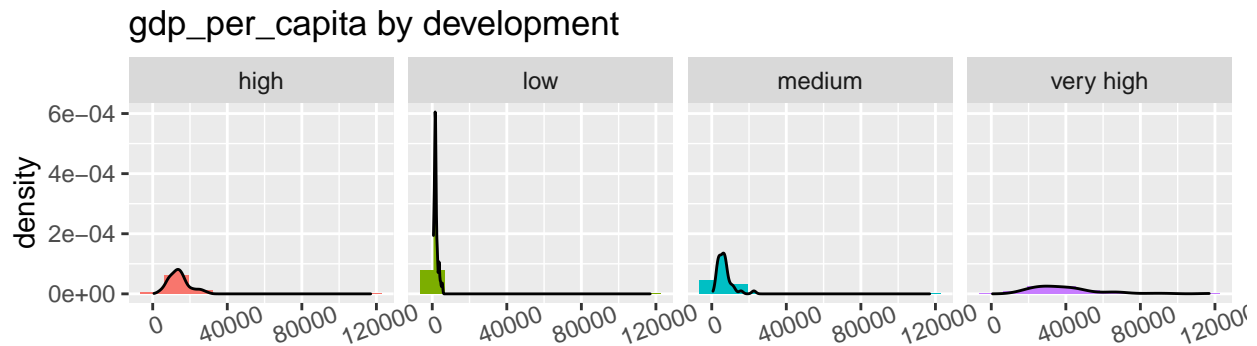
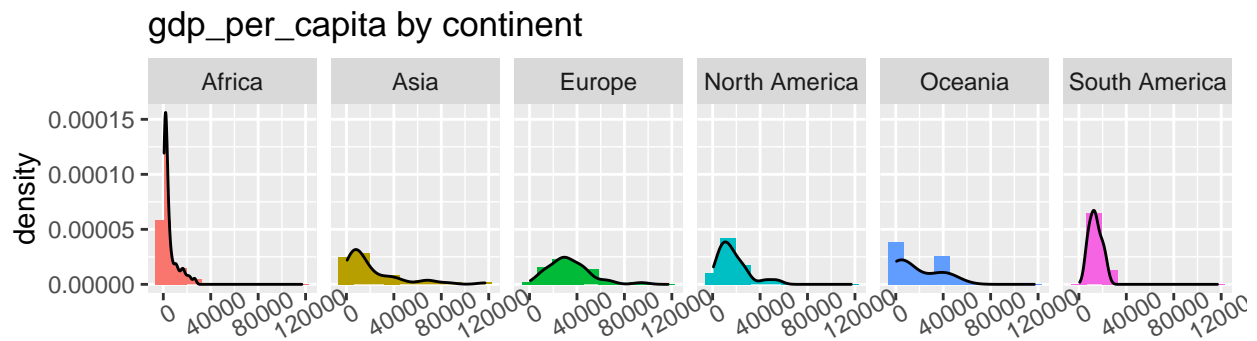
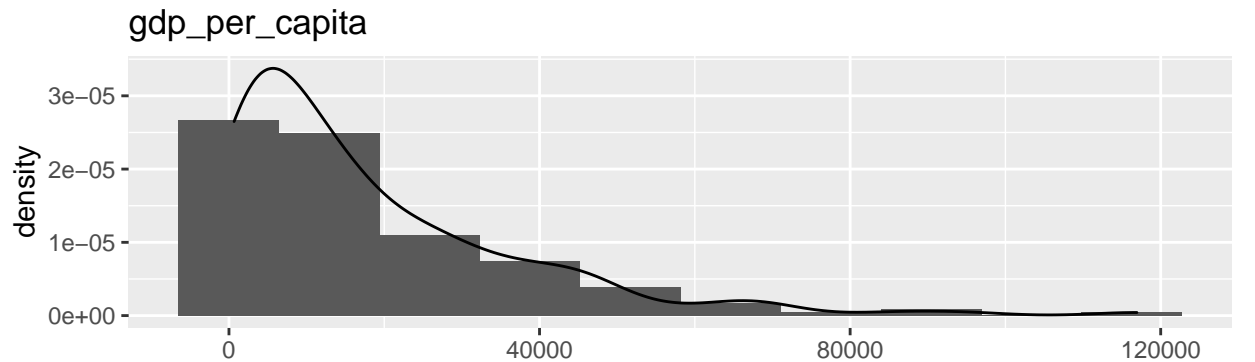
Observing the box-plot for the global GDP per capita we can see that the distribution is right skewed. The country that has the highest GDP per capita is Qatar from Asia, then comes Luxembourg and Singapore, all of them have a very high HDI.

The grouped box-plots provide more interesting conclusions: - In the box-plots grouped by continent we observe that the GDP per capita of Europe is a little bit higher than the rest of the continents while Africa has the least median of GDP per capita with some “outliers” that have similar values as other continents.

- Nevertheless, the box-plots grouped by development give us more relevant information. We can more or less define whether a new country has very high, high, medium or low HDI by having its GDP per capita. Due to the clear difference of GDP per capita between the different levels of HDI. The countries that have very high HDI often have larger GDP per capita, and the countries with low HDI have less GDP per capita. There is a very clear correlation between these two variables.

Histogram and kernel density for GDP per capita

```
plots(dataset=data, col='gdp_per_capita',type='hist', density=TRUE, xtick_angles=c(0,30,20))
```



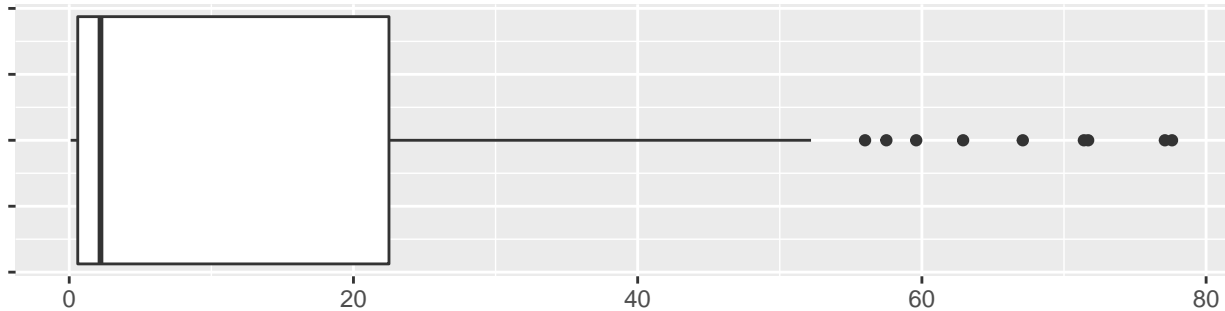
The global distribution of this variable is very right skewed. The distributions of Africa is more concentrated (leptokurtic) while others are more flat (platykurtic).

The distributions of countries that have low HDI are more concentrated, and they usually have less GDP per capita.

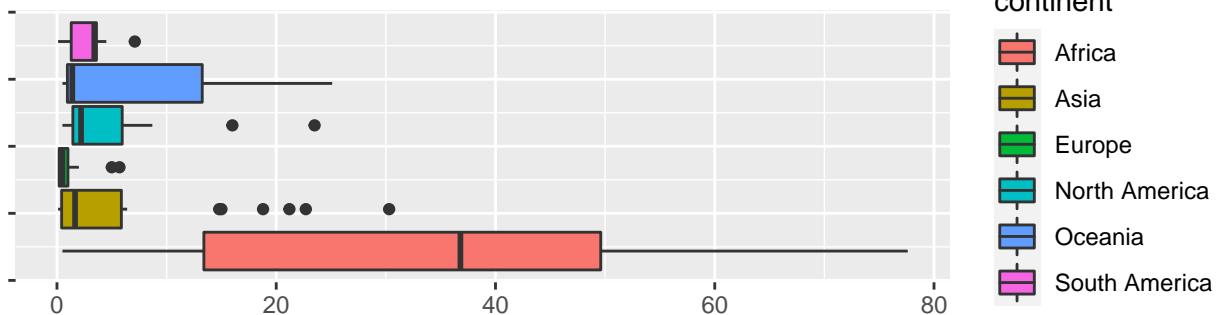
Boxplots for percentage of population in extreme poverty

```
plots(dataset=data, col='extreme_poverty',type='boxplot')
```

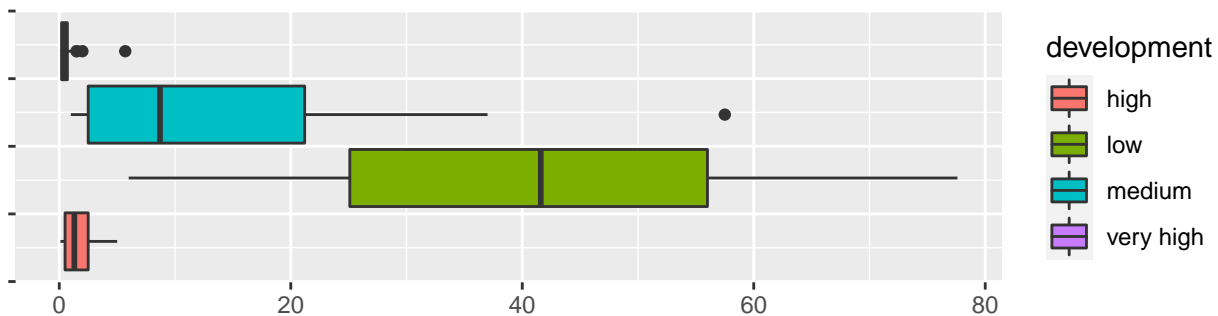
extreme_poverty



extreme_poverty grouped by continent



extreme_poverty grouped by development



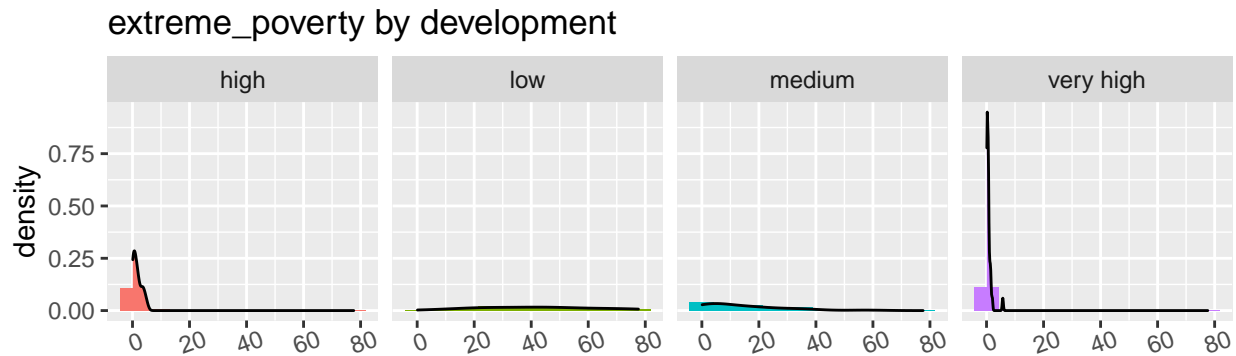
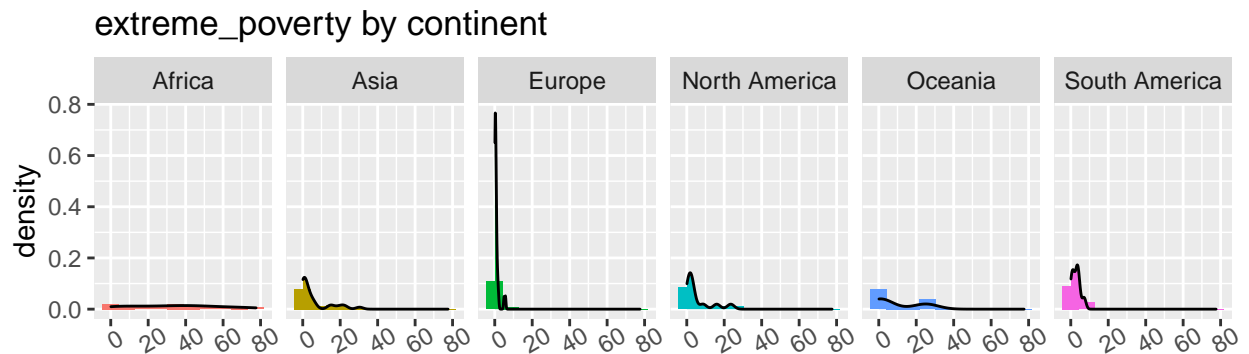
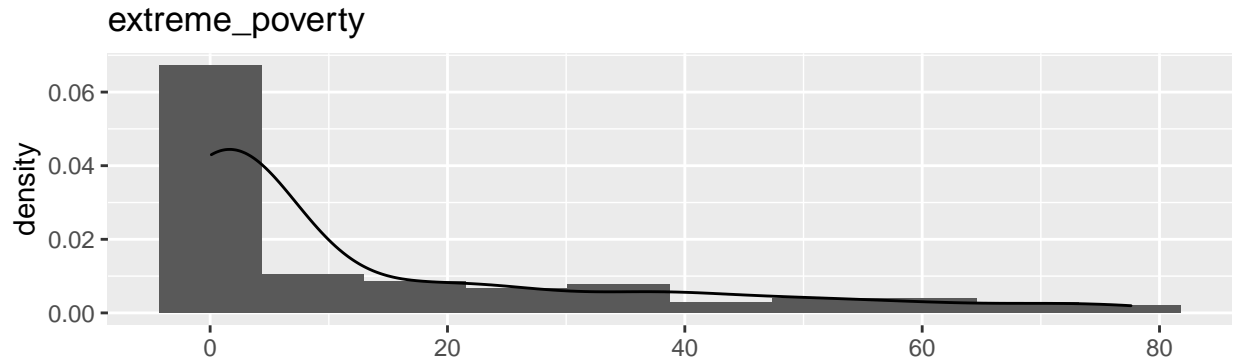
From the box-plot of the global extreme poverty we can observe that the distribution is right skewed. The country that has the highest extreme poverty is Madagascar, then comes Democratic Republic of Congo and Burundi, all of them are from Africa.

However, the grouped box-plots provide more interesting conclusions: - The box-plots grouped by continent tell us that the median of the extreme poverty of Europe is the least of all the continents while Africa has the highest median extreme poverty.

- The box-plots grouped by development give us more important information. The countries that have higher index of extreme poverty often have low HDI while the countries with lower extreme poverty have higher HDI. There is a quite clear correlation between these two variables.

Histogram and kernel density for percentage of population in extreme poverty

```
plots(dataset=data, col='extreme_poverty',type='hist', density=TRUE, xtick_angles=c(0,30,20))
```

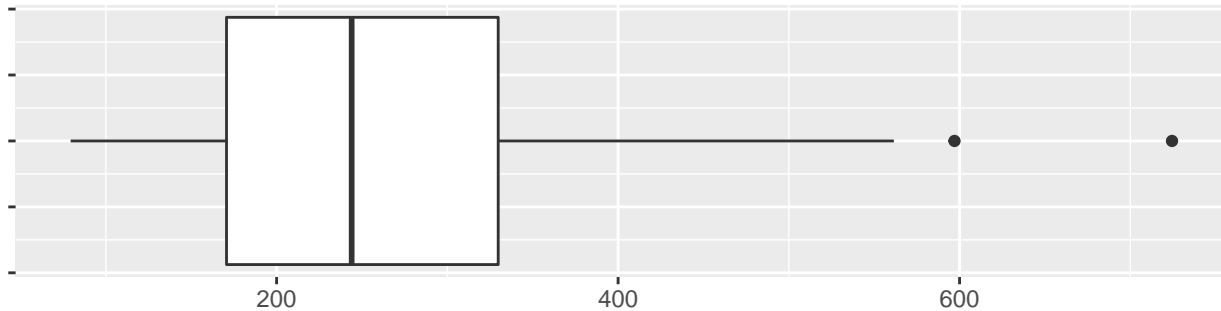


The global distribution of this variable is very right skewed. The distributions of Europe is more concentrated (leptokurtic) in low values while others are more flat (platykurtic). The distributions of countries that have very high HDI are more concentrated, and they usually have lower extreme poverty.

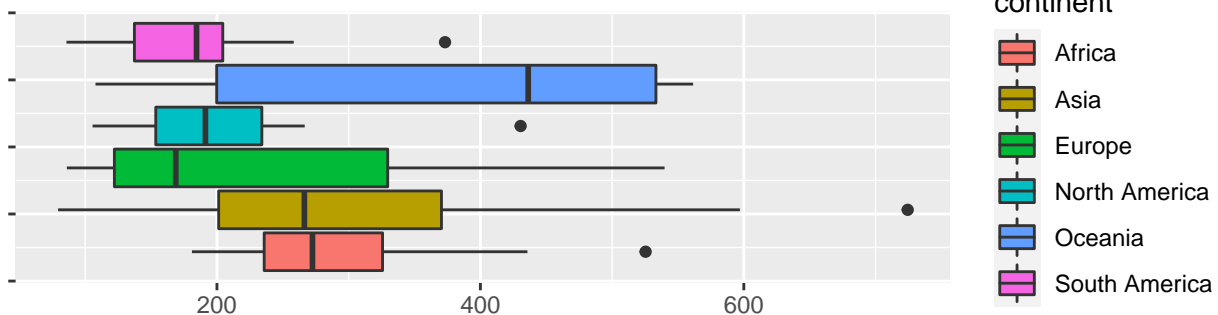
Boxplots for cardiovascular death rate

```
plots(dataset=data, col='cardiovasc_death_rate',type='boxplot')
```

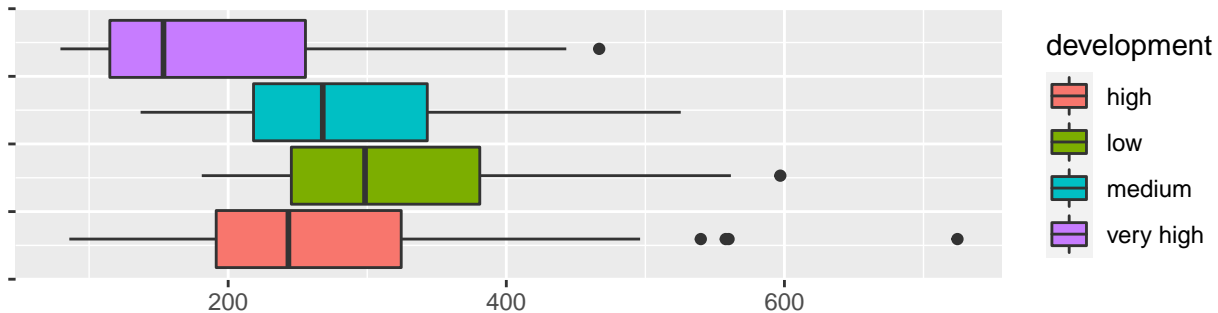
cardiovasc_death_rate



cardiovasc_death_rate grouped by continent



cardiovasc_death_rate grouped by development



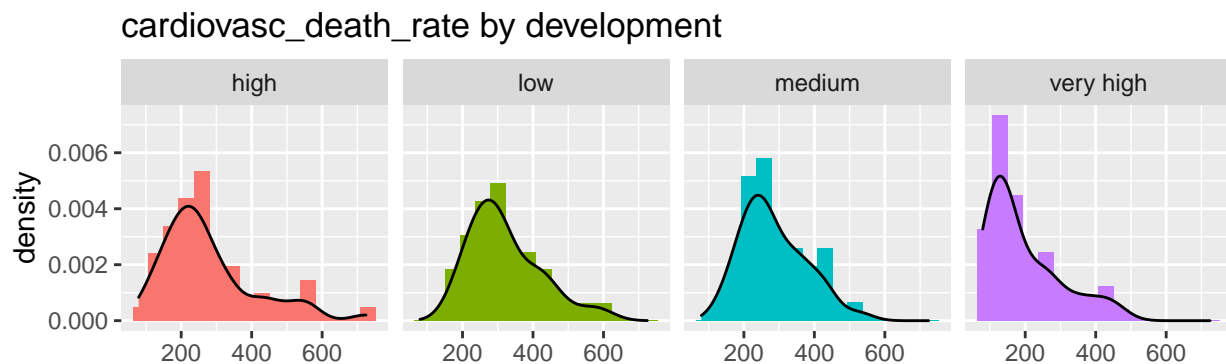
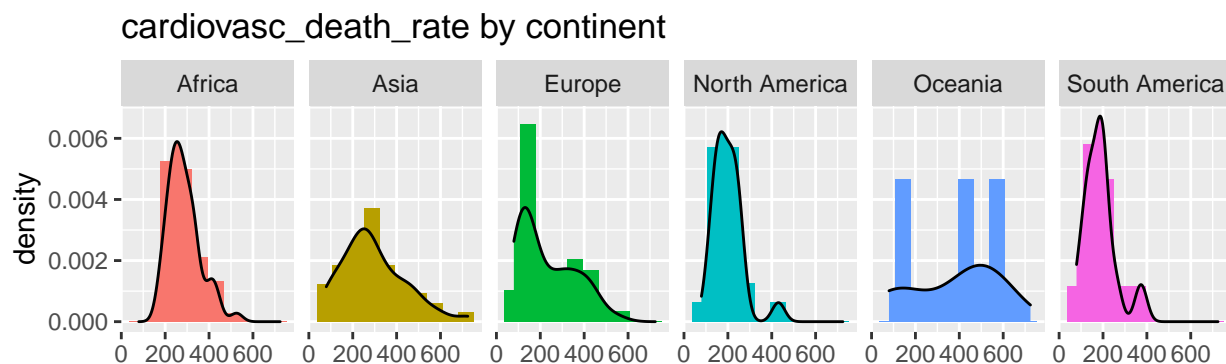
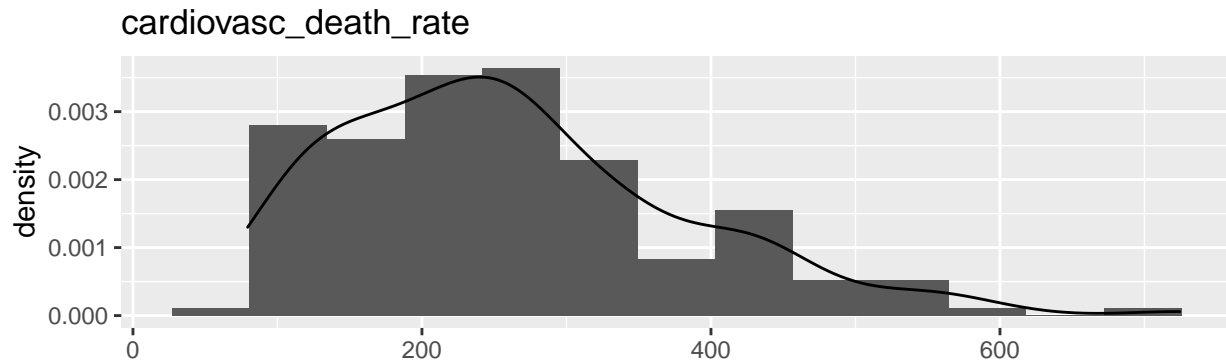
It appears that most countries seem to have a cardiovascular death rate between 170.67 and 329.79 deaths per 100,000 inhabitants. With Uzbekistan being in the absolute extreme, with about 724 deaths by cardiovascular disease per 100,000 inhabitants.

Grouping by continent we see that Oceania seems to have the largest box (probably due to its lower amount of countries), with a few extreme cases per continent. On average the continent with the highest death rate due to cardiovascular disease is Oceania, followed by Asia. Seems like cardiovascular disease in the Americas could be a less common cause of death than in the rest of the world.

By development we can see a bit of a pattern, where the least developed a country is, the higher its cardiovascular death rate. However, even if we see this pattern, we can't confidently say that living in a less developed country makes an individual more likely to die from cardiovascular disease. There are definitely many other factors that affect such rate per HDI.

Histogram and kernel density for cardiovascular death rate

```
plots(dataset=data, col='cardiovasc_death_rate', type='hist', density=TRUE, bins=c(13,10,16), xtick_angle=45)
```



For cardiovascular death rate we see a similar story here than with the general boxplot. The larger concentration of countries clumps around the previously mentioned interval, and the distribution of the variable as is is somewhat normal-like with a relatively long left tail.

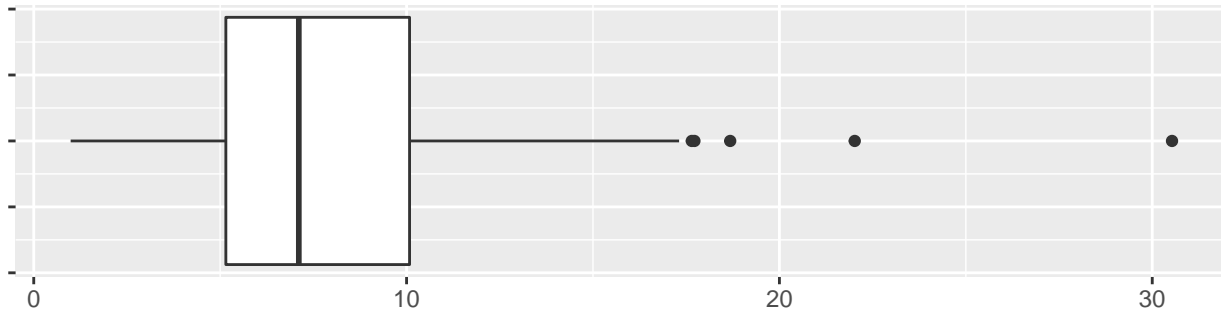
Per continent we see that Europe has a significant concentration of countries below 200, along with South America, which, on average, is the continent with the lowest death rate from cardiovascular disease according to our data. For Oceania we see a flat distribution with some high numbers and low numbers, of course, we know that there's less data points, therefore our main concentration below 200 corresponds to New Zealand and Australia, and the rest of the countries seem to have a higher death rate than the rest. Asia's left tail suggests a few other countries with a very high cardiovascular death rate like Uzbekistan.

Looking at development we see the much higher concentration of low cardiovascular death rates for very high development countries. Which in general tend to have better healthcare. However while lower for low and medium development countries, we don't see too much of a difference between the two in terms of their distribution.

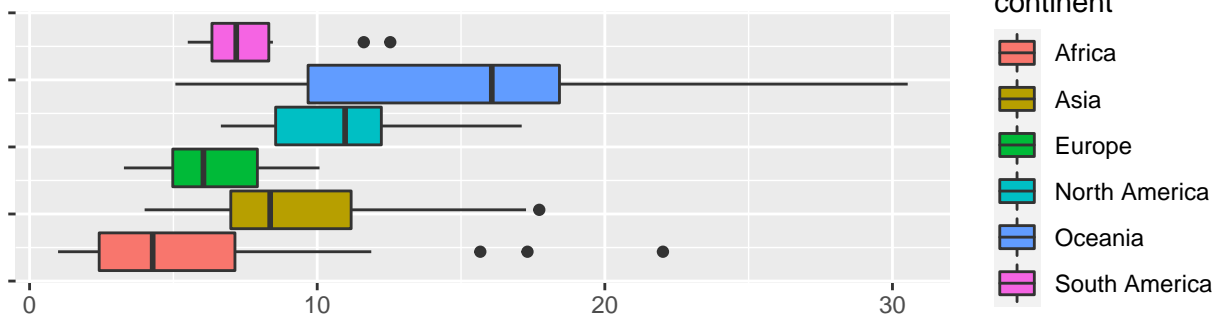
Boxplots for diabetes prevalence

```
plots(dataset=data, col='diabetes_prevalence',type='boxplot')
```

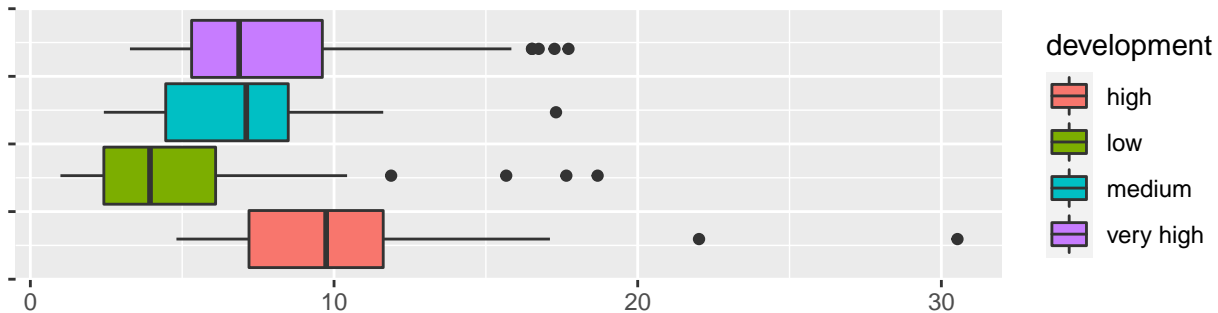
diabetes_prevalence



diabetes_prevalence grouped by continent



diabetes_prevalence grouped by development



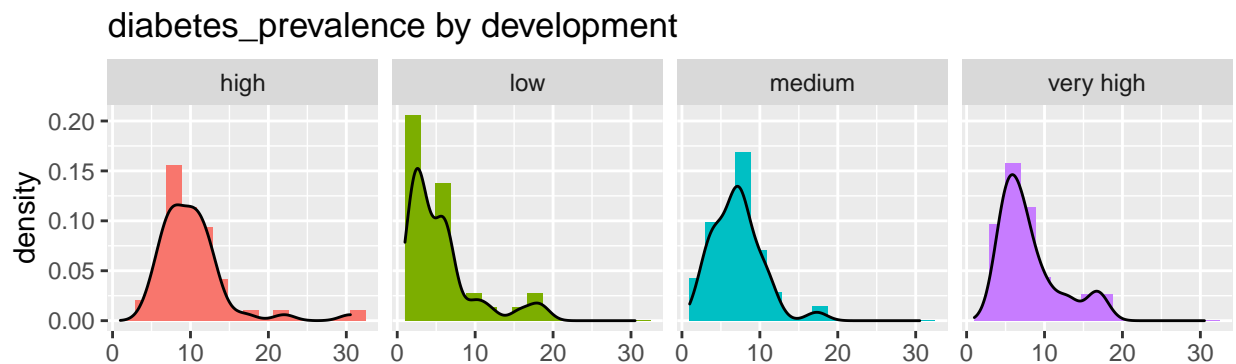
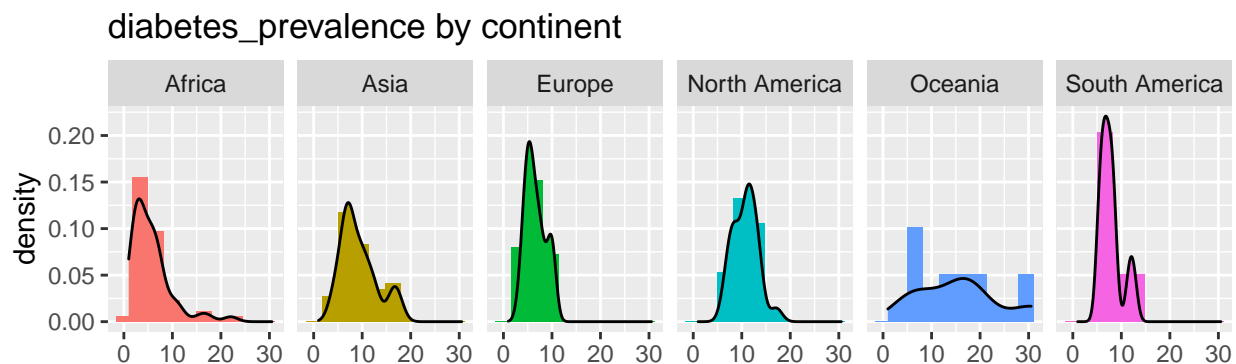
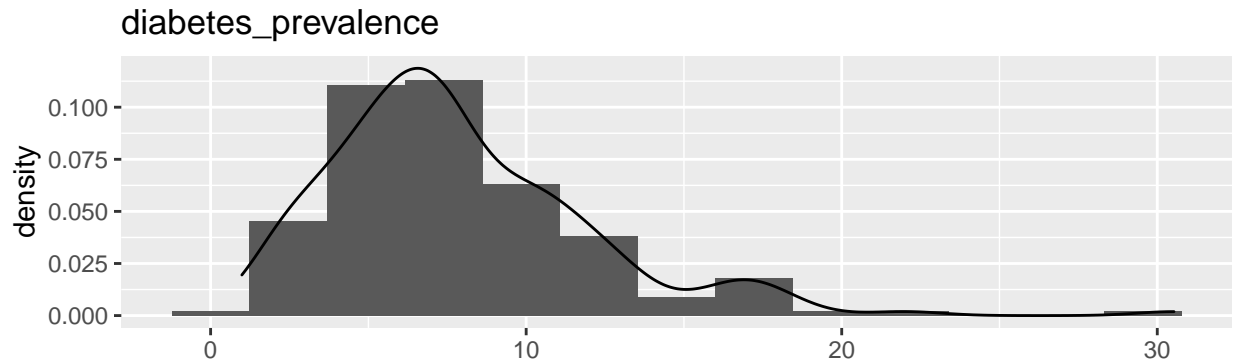
For diabetes prevalence we can see that most countries sit at a value of less than 10, but higher than 5. Some countries surpassing even 30%. These extreme values correspond to a few countries in Oceania and Africa.

The continent with the highest incidence as a proportion of its population seems to be Oceania. Which includes the top 3 countries with the highest amount of diabetics as a percentage of their population. With 30.53% for Marshall Islands. Although the values of the other top 2 countries are not included in our dataset, after some research, we found out that they're also 2 countries in Oceania. North America's diabetes incidence has nearly doubled in the past 20 years, therefore taking the spot 2 as the continent with the highest incidence with Asia, South America, Europe and Africa trailing behind.

We can, to an extent, see that higher development doesn't necessarily mean higher or lower diabetes prevalence and this might relate more to genetic composition and diet of the inhabitants.

Histogram and kernel density for diabetes prevalence

```
plots(dataset=data, col='diabetes_prevalence', type='hist', density=TRUE, bins=c(13,10,16), xtick_angles=
```



For the distribution of the data we see that it resembles a normal distribution with a long left tail and the most countries clumped around the mean of ~7.9%.

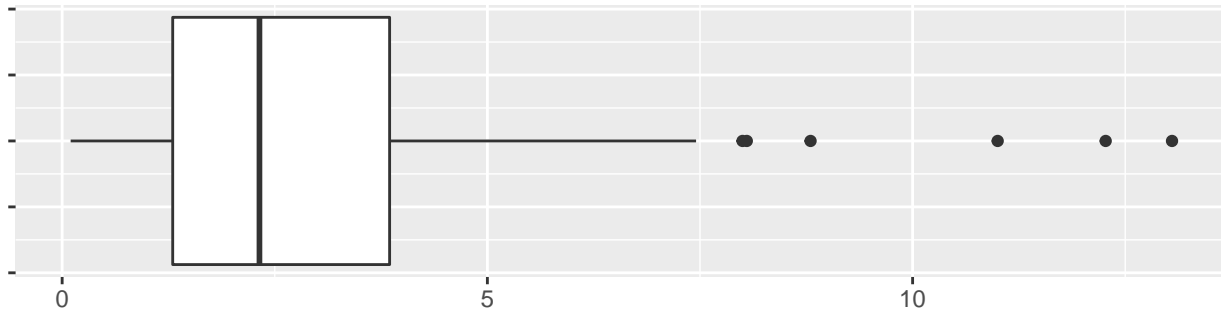
For each continent the incidence seems to be quite different, with some continents having a much higher incidence than others (for example Oceania vs Africa), however they all seem to clump around similar values.

For the development we see the same we saw in the boxplots. Not much of a pattern or indication that there's any specific relationship between HDI and diabetes prevalence.

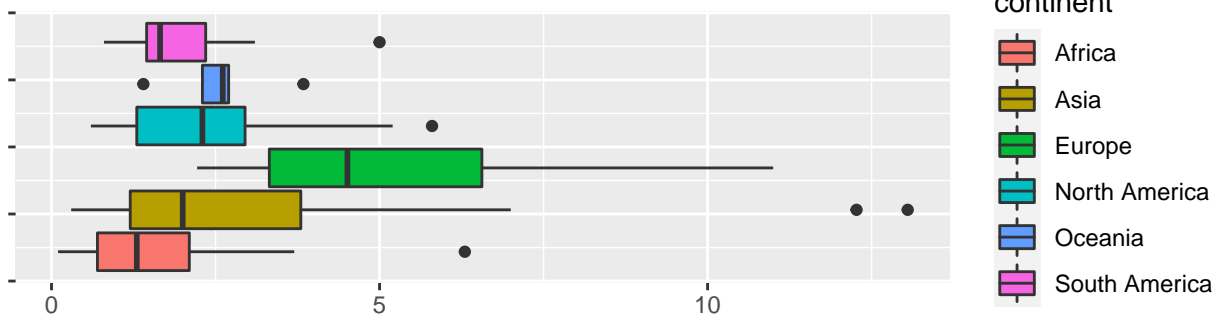
Boxplots for hospital beds per thousand inhabitants

```
plots(dataset=data, col='hospital_beds_per_thousand',type='boxplot')
```

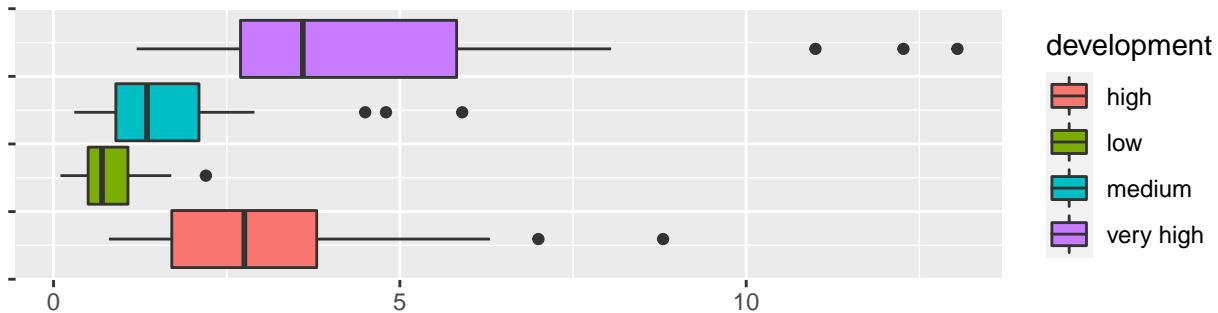
hospital_beds_per_thousand



hospital_beds_per_thousand grouped by continent



hospital_beds_per_thousand grouped by development



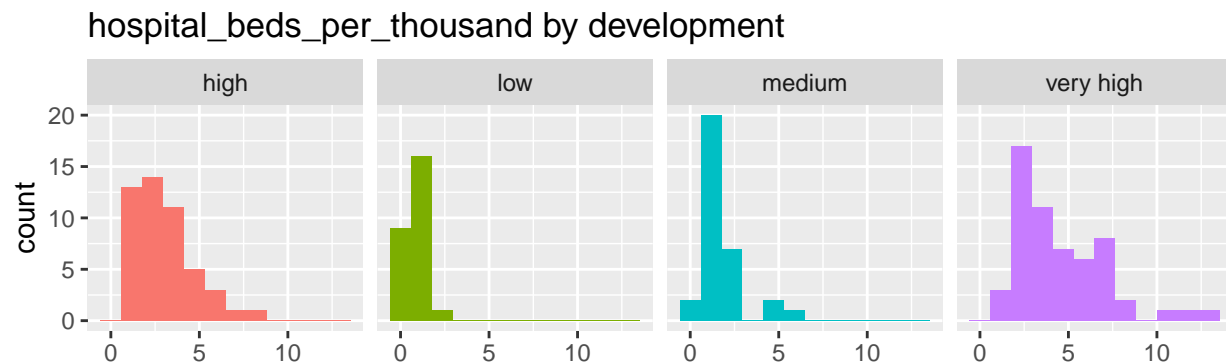
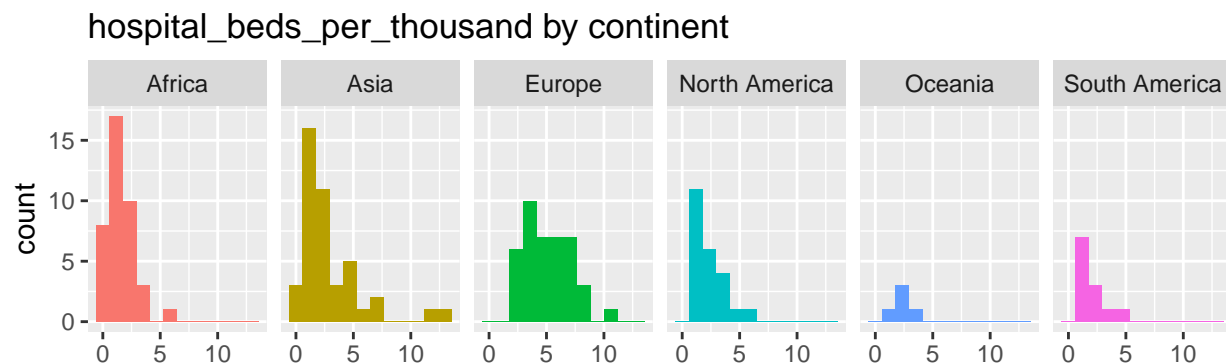
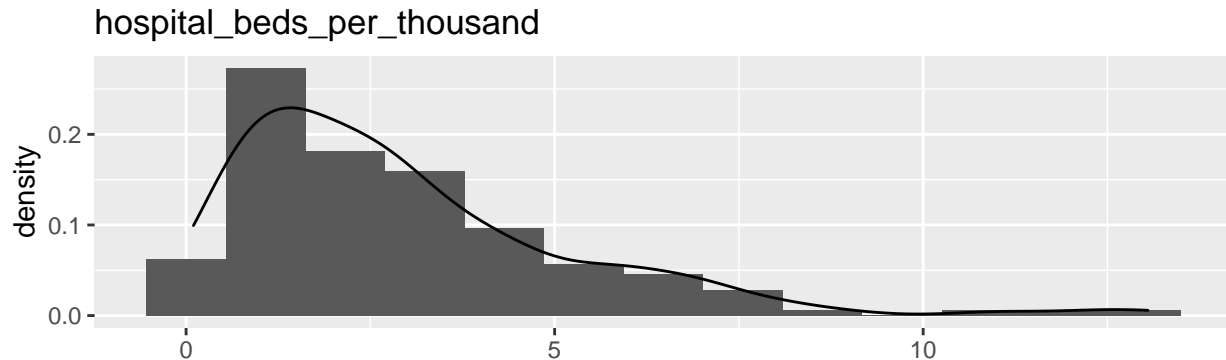
Looking at the hospital beds per thousand inhabitants variable boxplots we can see a few interesting things. We could use this variable as a measure of the quality of a healthcare system of a country. Where the higher the bed availability in hospitals is, the better the health system can cope with the demand for beds that a pandemic usually comes with. Especially with how widespread COVID-19 is.

We can see that some extremely poor countries have about 0.1-0.3 beds per thousand inhabitants, like it is the case with Mali and Niger. Some other countries like South Korea or Belarus have an extremely high capacity, with around 12 and 11 beds per thousand inhabitants respectively. However, even if the amount of beds per thousand inhabitants seems to be low, there's some countries with a suspicious seemingly low amount of beds, however, some of this are clearly just very highly populated countries.

For countries with high and very high HDI, there's a clear bias towards having greater bed capacity, however, this is not the case for all countries with that quality as there's clearly some countries with medium HDI that have a quite formidable bed capacity as well.

Histogram and kernel density for hospital beds per thousand inhabitants

```
plots(dataset=data, col='hospital_beds_per_thousand',type='hist', density=FALSE, bins=c(13,12,12), xticl
```



These plots tell a little bit of a different story to the boxplots. Where the largest concentration of countries is between 0 and 5 hospital beds per thousand inhabitants with an extremely scarce amount of countries with more than 10 beds per thousand inhabitants.

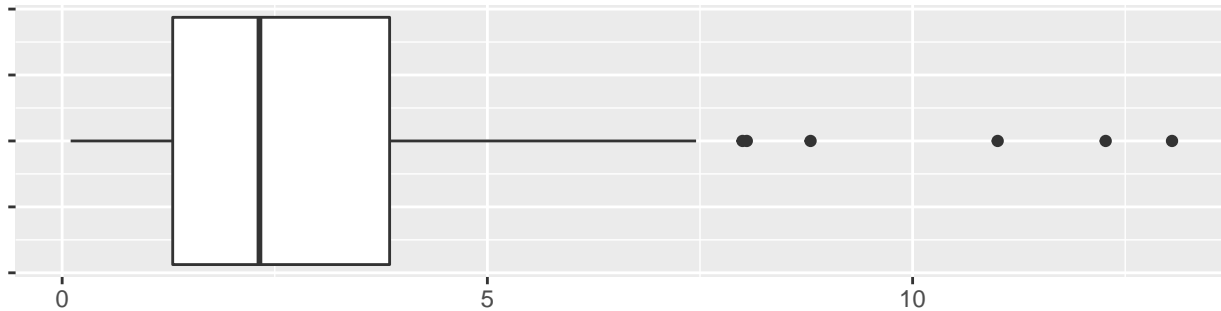
By segregating the data by continent we see that development does not necessarily mean greater healthcare capacity, with most continents boasting very similar numbers in this aspect while some like Asia, Africa, Europe and North America possessing some exceptions with extremely high numbers compared to the rest. However, yes, there's definitely a hint in continents with more developed countries (like Europe or some parts of Asia) which have a higher amount of beds, while Africa, which is predominantly composed of less developed countries tend to have a lower amount of beds.

Finally, looking at development we see that it is rare for much less developed countries to have high bed capacity, while it is much easier for high to very high developed countries to have greater capacity. However, we can't confidently say that there's lots of exceptions to this 'rule'.

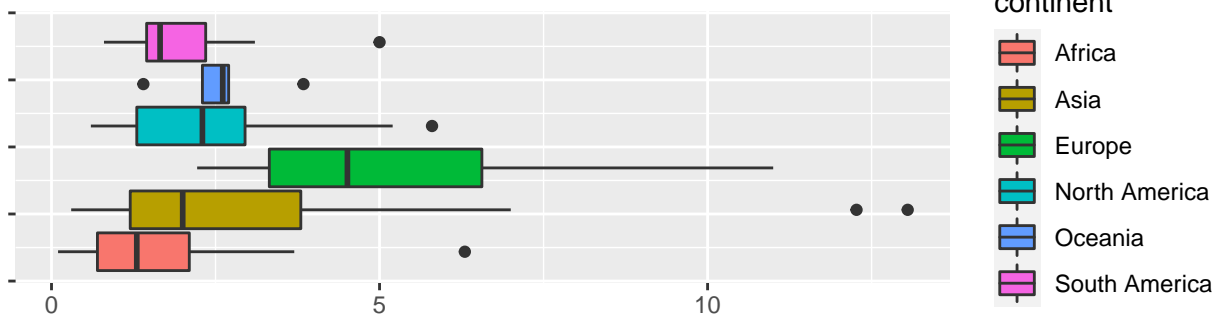
Boxplots for life expectancy

```
plots(dataset=data, col='hospital_beds_per_thousand',type='boxplot')
```

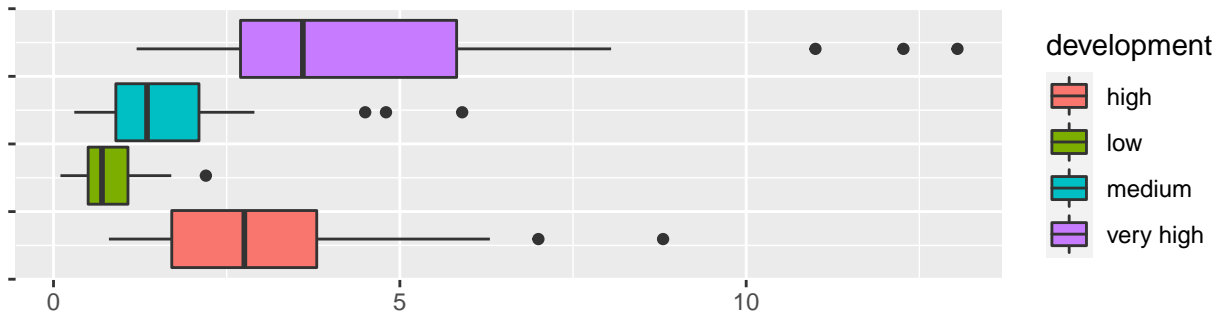
hospital_beds_per_thousand



hospital_beds_per_thousand grouped by continent



hospital_beds_per_thousand grouped by development



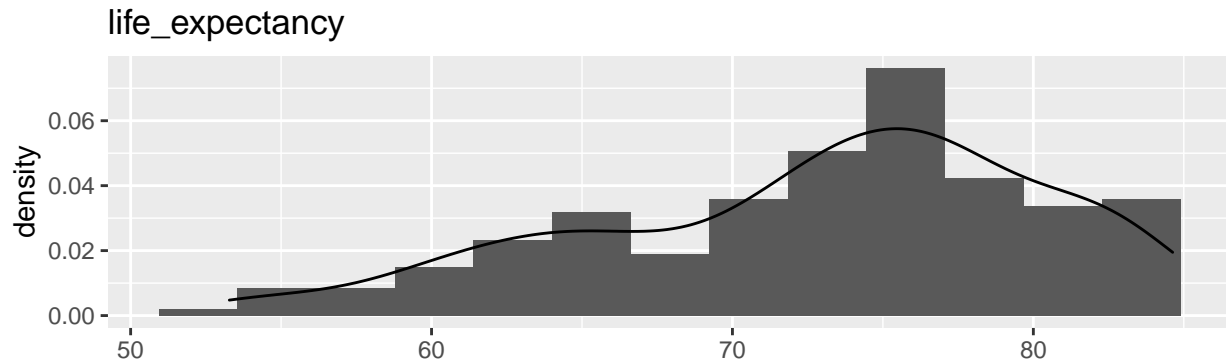
For life expectancy we can see most countries sitting above 66 years of age, with values going as low as 53.24 and as high as 84.63.

Africa has the lowest life expectancy while Europe has the highest. The rest of the continents sit at roughly similar ranges.

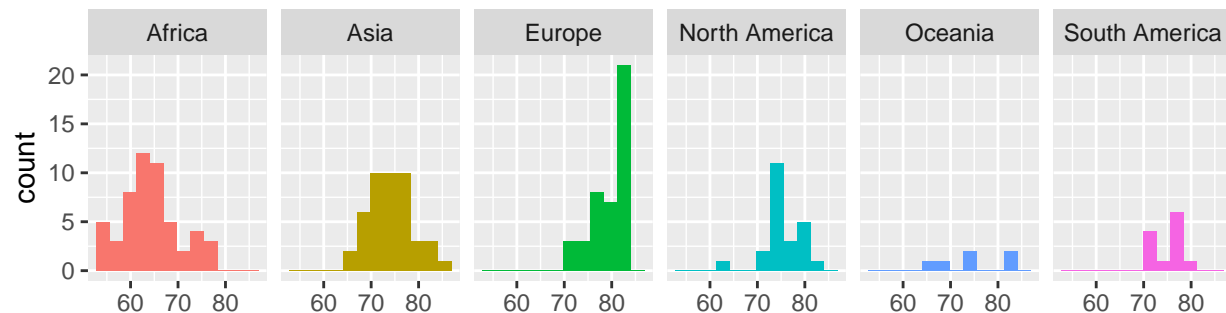
Grouping by HDI, we can see that the most developed countries have a significantly higher life expectancy than those with low HDI. It clearly shows a strong positive correlation between them. Where the higher the life expectancy the higher the HDI. With very few exceptions.

Histogram and kernel density for life expectancy

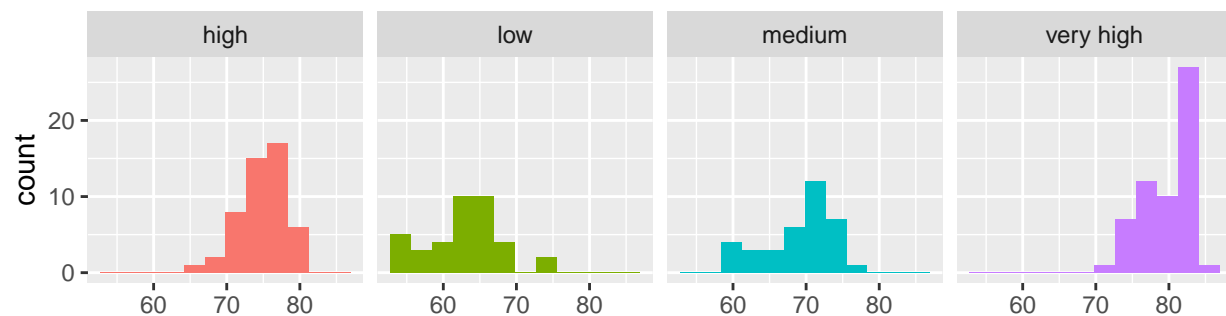
```
plots(dataset=data, col='life_expectancy',type='hist', density=FALSE, bins=c(13,12,12), xtick_angles=c(
```



life_expectancy by continent



life_expectancy by development



The general plot is somewhat left skewed, as most countries (about 80%) have a life expectancy higher than 65 years of age. Our density plot shows a strong concentration between 70 and 80 years of age, as this range covers the most nations.

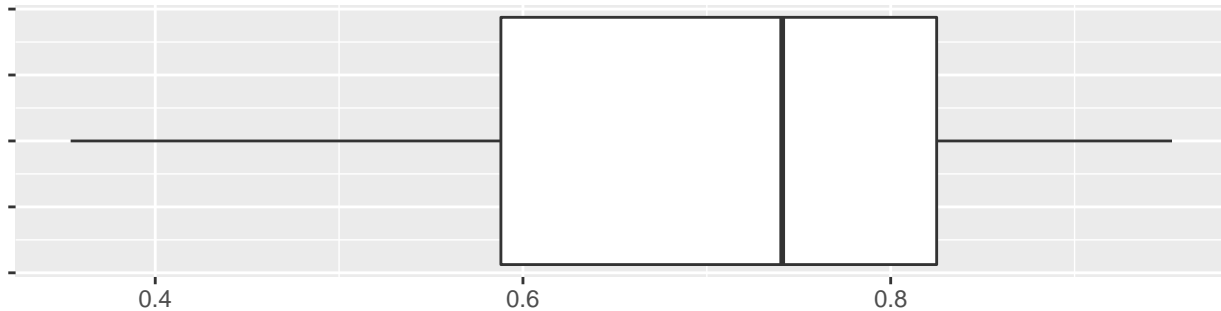
For each continent we see that Europe shows a typically very high life expectancy while Africa shows a typically lower-than-average life expectancy for most countries with some exceptions. The rest of the continents sit at about the average life expectancy with some countries in Asia and North America at significantly higher-than-average numbers.

For HDI we can again see some of the strong correlation, where life expectancy for very highly developed nations seems to be also quite high and the same happens with less developed nations.

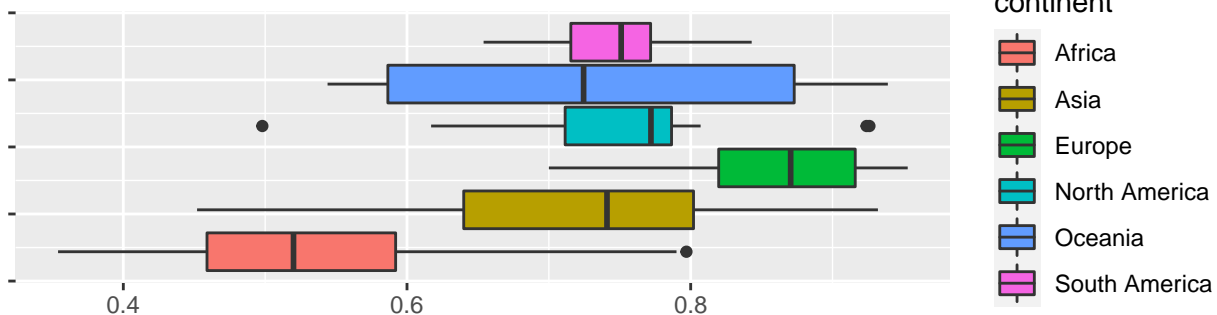
Boxplots for Human Development Index

```
plots(dataset=data, col='human_development_index', type='boxplot')
```

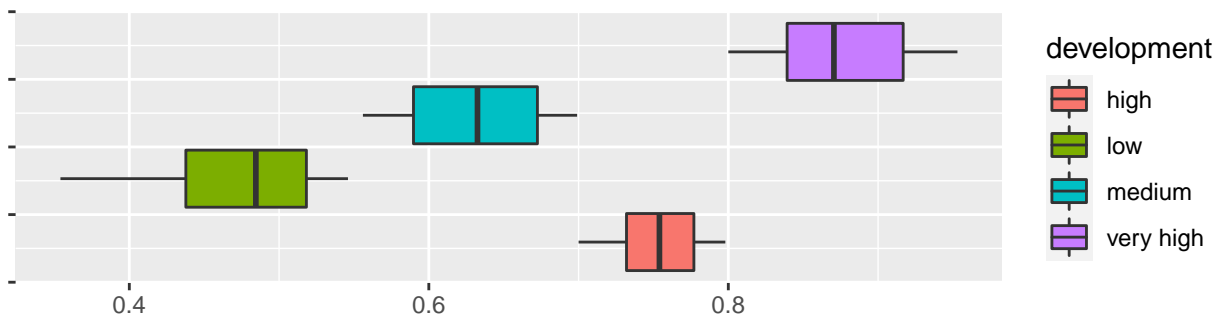
human_development_index



human_development_index grouped by continent



human_development_index grouped by development



We can see most countries fall between 0.6 and 0.8, our median HDI is 0.741.

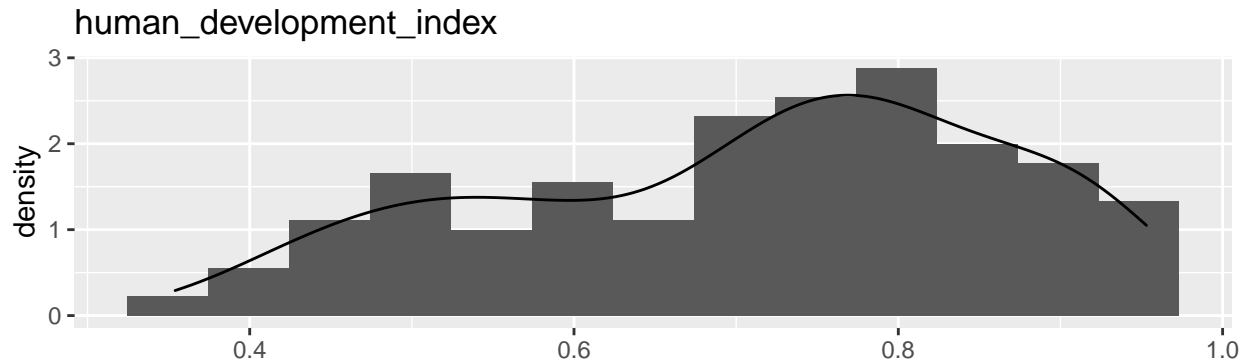
For continents we can see Africa lagging behind with most of its countries between 0.4 and 0.6 HDI, probably given the poverty situation in the continent.

The rest of the continents sit between 0.6 and 0.8 for most of its countries with North America having 2 very extreme outliers which are its minimum and maximum values (corresponding respectively to Haiti and USA). Europe is generally above 0.8.

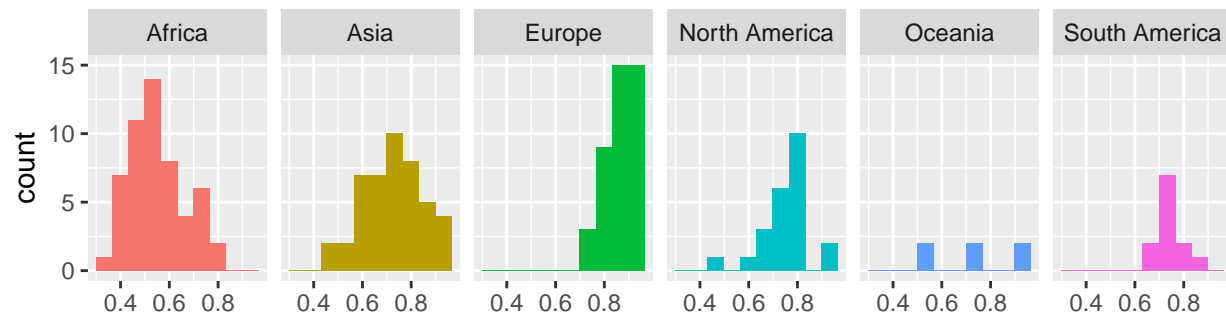
As our development variable was constructed from the human_development_index variable, we can see that there's clearly marked bounds for each HDI range. The ranges are as follows: *very high* for HDI of 0.800 and above, *high* from 0.700 to 0.799, *medium* from 0.550 to 0.699 and *low* below 0.550.

Histogram and kernel density for Human Development Index

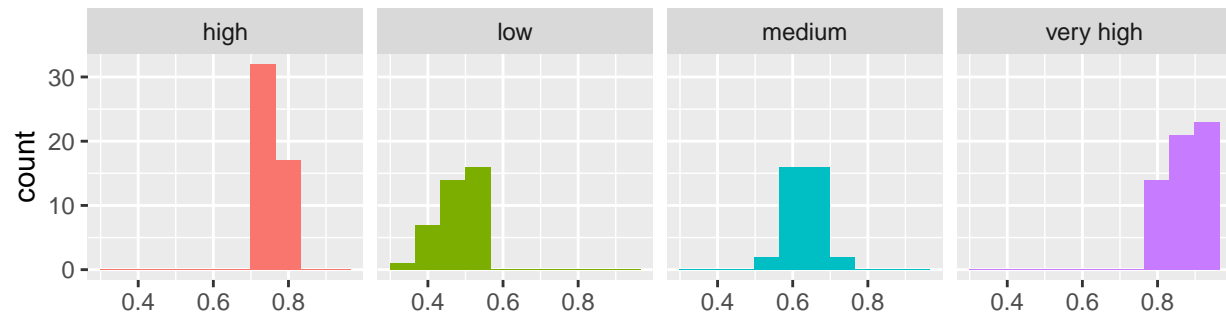
```
plots(dataset=data, col='human_development_index', type='hist', density=FALSE, bins=c(13,10,10), xtick_a
```



human_development_index by continent



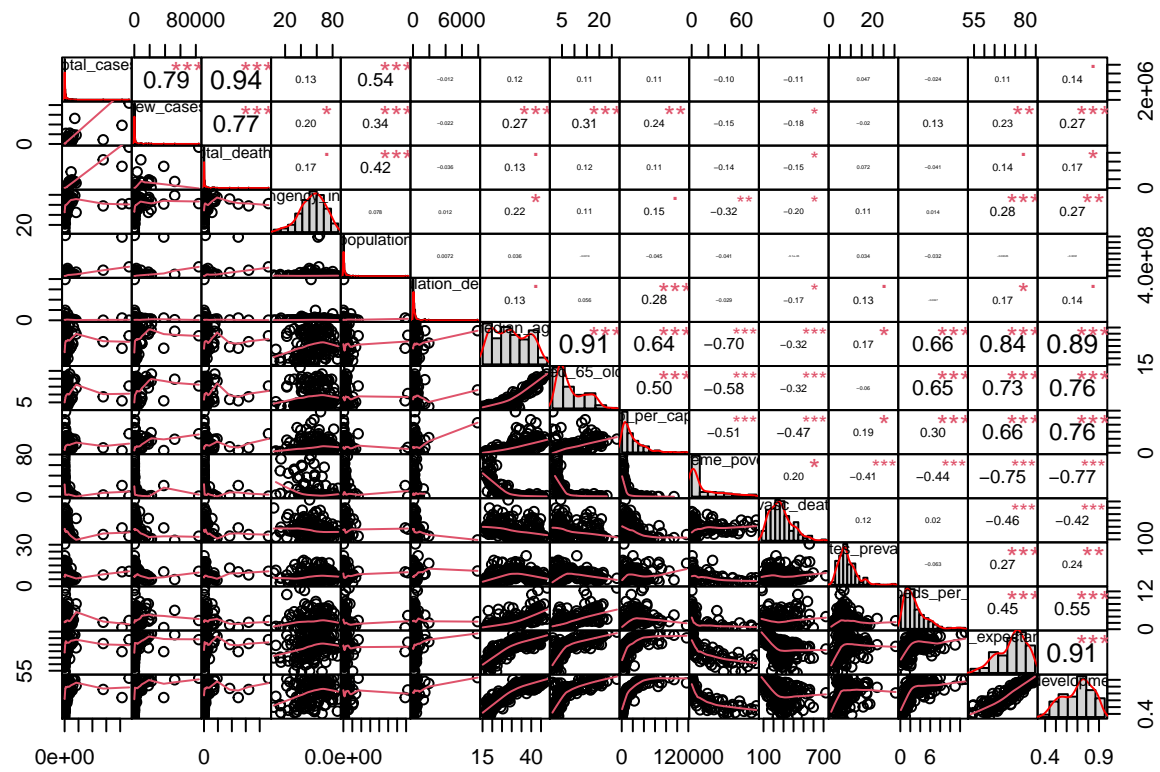
human_development_index by development



For the human development index we can see that the variable is somewhat left skewed, given that the average HDI is ~ 0.71 , which most countries either match or are above of.

For the HDI per continent we can see that africa has a clear concentration below 0.6, given that most countries in Africa have a low HDI. South america and Asia tell a similar story, most countries are at or above 0.6. We can see that for North America there's a little concentration below 0.6 and most countries between 0.6 and 0.8 as North America includes Central America and the Caribbean which tend to have a lower HDI than USA/Canada, which are towards the right of 0.8. Most european countries have a very high to high HDI, therefore the density plot is quite left skewed and most contries in Oceania have a lower-than-average HDI with the exception of New Zealand and Australia which are above 0.8.

```
pa <- data_n %>% dplyr::select(interesting_vars)
chart.Correlation(pa, histogram=TRUE, pch=19, method="pearson")
```



PCP Plot

We define a function to set colors for categorical variables in a PCP plot:

```
colors <- function(cat_var, colors_vector) {  
  kleuren <- as.numeric(as.factor(cat_var))  
  foreach (i=1:length(kleuren), kleur=kleuren) %do% {  
    kleuren[i] = colors_vector[kleur]  
  }  
  return(kleuren)  
}
```

Colours we picked:

```
# setting colors development  
color_1 <- "blueviolet"  
color_2 <- "brown"  
color_3 <- "seagreen"  
color_4 <- "yellow3"  
color_5 <- "black"  
color_6 <- "deeppink1"  
palette1 <- c(color_1,color_2,color_3,color_4)  
palette2 <- c(color_1,color_2,color_3,color_4,color_5,color_6)
```

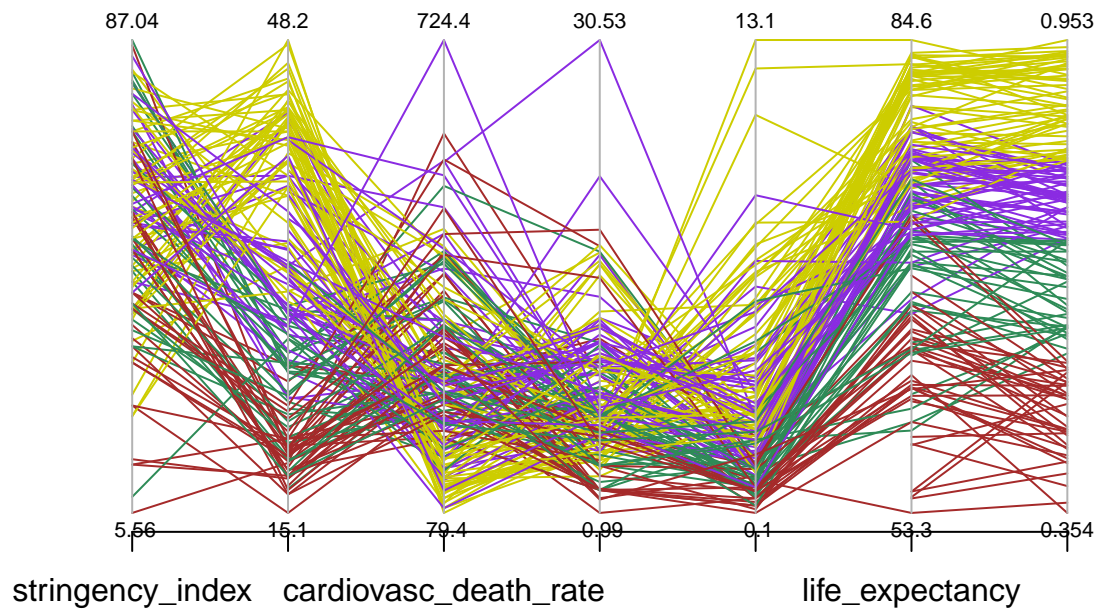
```
development_colors <- colors(data$development,palette1)  
continent_colors <- colors(data$continent,palette2)
```

We group variables by their skewness, while we have many right skewed variables, we group the rest of them in another PCP plot, to have a less crowded plot.

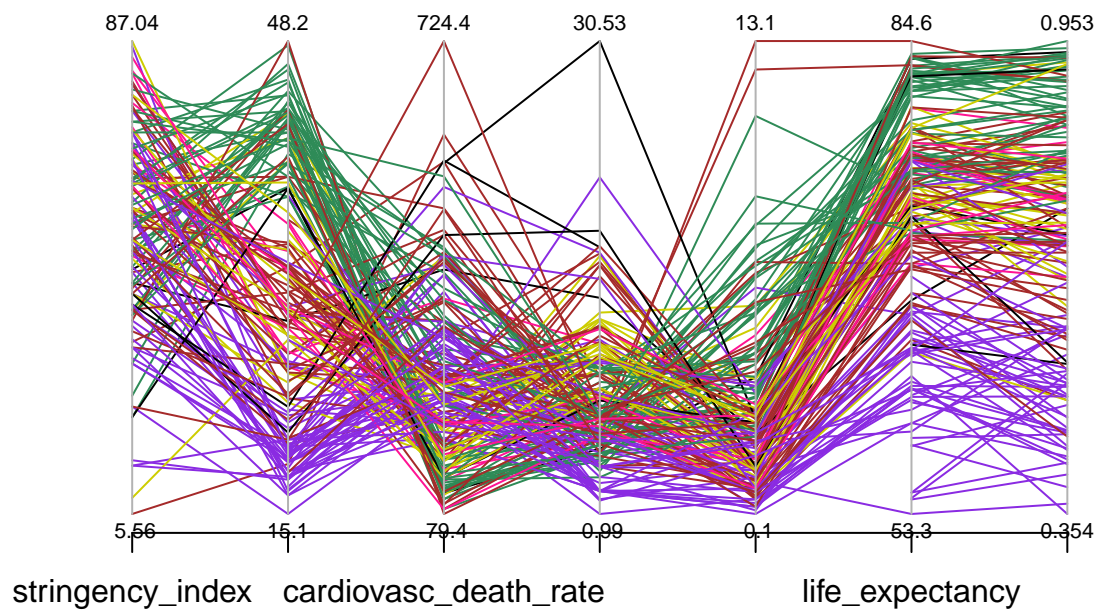
```
right_skewed <- c('total_cases','new_cases','population',  
                 'population_density','aged_65_older',  
                 'gdp_per_capita','extreme_poverty')  
right_skewed <- data_n %>% dplyr::select(right_skewed)  
  
others <- c('stringency_index','median_age',  
           'cardiovasc_death_rate','diabetes_prevalence',  
           'hospital_beds_per_thousand','life_expectancy',  
           'human_development_index')  
others <- data_n %>% dplyr::select(others)
```

Right skewed variables PCP

```
parcoord(others,var.label=TRUE, col=development_colors)
```

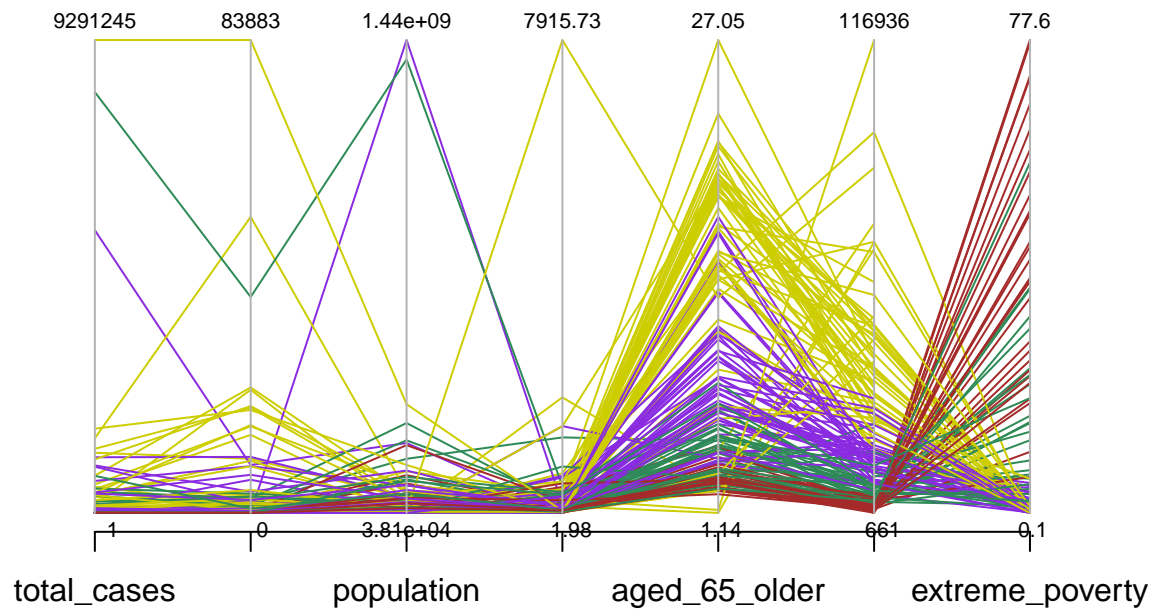


```
parcoord(others,var.label=TRUE, col=continent_colors)
```

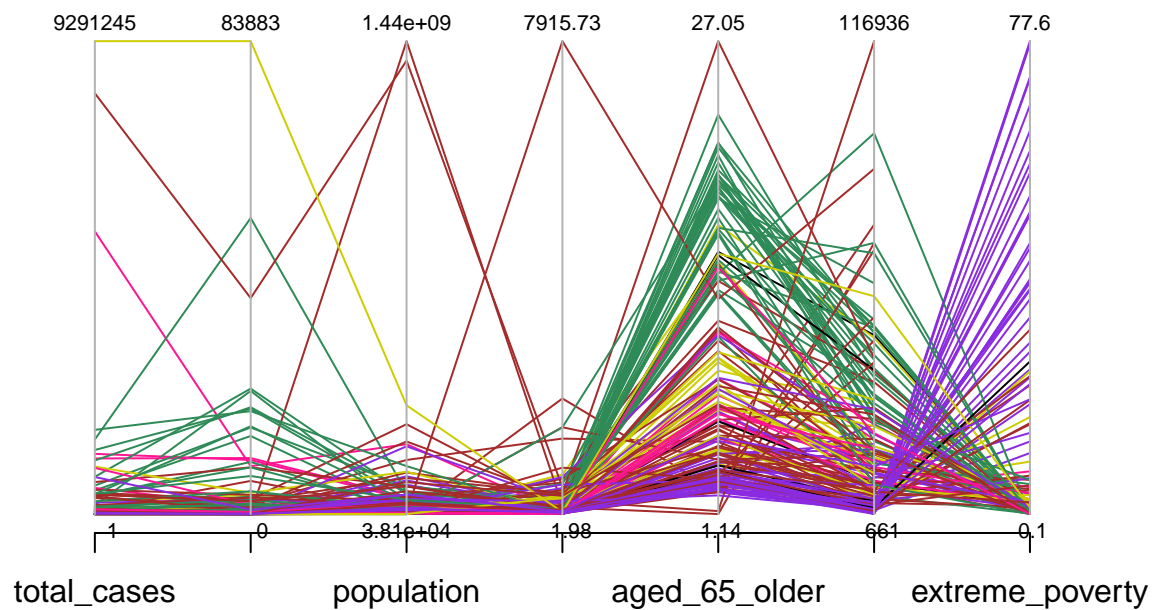


Other variables PCP

```
parcoord(right_skewed,var.label=TRUE, col=development_colors)
```



```
parcoord(right_skewed,var.label=TRUE, col=continent_colors)
```



```
library(mice)
covid = data[c('continent','location','development',
               'total_cases','new_cases','total_deaths',
               'stringency_index','population',
               'population_density','median_age',
               'aged_65_older','gdp_per_capita',
               'extreme_poverty','cardiovasc_death_rate',
               'diabetes_prevalence','hospital_beds_per_thousand',
               'life_expectancy','human_development_index')]
covid$continent=factor(covid$continent)
covid$development=factor(covid$development)
covid_imp=mice(covid,m=5,method = "cart")
covid_imp=complete(covid_imp)
```