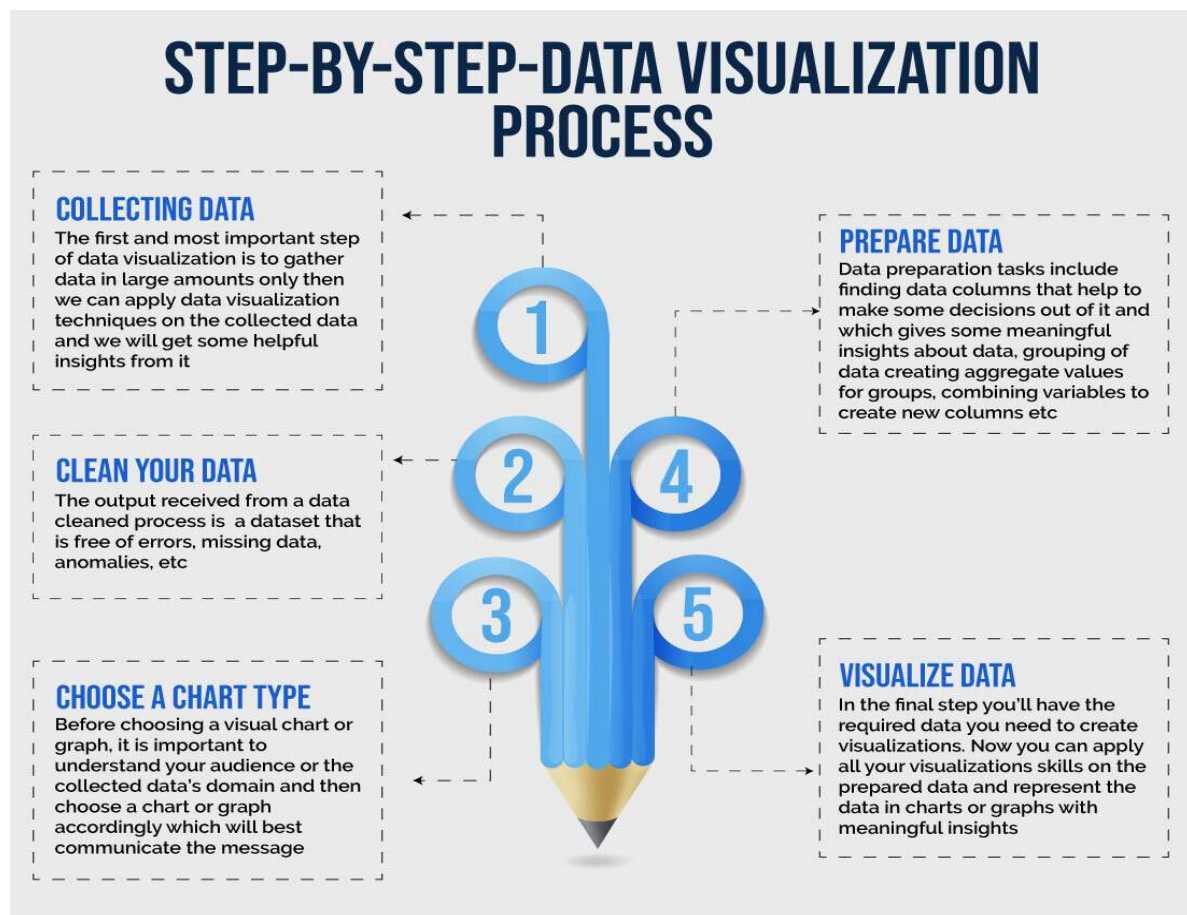# DATA VISUALISATION

## Introduction:

Three factors define data visualisation. The data visualisation method generates an image that is legible by viewers and helps the investigation, examination, and transmission of the information based on either qualitative or quantitative data (Azzam et al,2013).

By adhering to data visualisation guidelines, a scientist can enhance their visual message. Other principles are more technical, such as what colour combinations transmit information, while others are less technical, such as choosing a message before starting the graphic. Due to their lack of formal training in figure creation, scientists must be aware of best practises in order to tell the story of the data they

## STEP-BY-STEP-DATA VISUALIZATION PROCESS

**COLLECTING DATA**
The first and most important step of data visualization is to gather data in large amounts only then we can apply data visualization techniques on the collected data and we will get some helpful insights from it

1

**CLEAN YOUR DATA**
The output received from a data cleaned process is a dataset that is free of errors, missing data, anomalies, etc

2

4

3

5

**CHOOSE A CHART TYPE**
Before choosing a visual chart or graph, it is important to understand your audience or the collected data's domain and then choose a chart or graph accordingly which will best communicate the message

**PREPARE DATA**
Data preparation tasks include finding data columns that help to make some decisions out of it and which gives some meaningful insights about data, grouping of data creating aggregate values for groups, combining variables to create new columns etc

**VISUALIZE DATA**
In the final step you'll have the required data you need to create visualizations. Now you can apply all your visualizations skills on the prepared data and represent the data in charts or graphs with meaningful insights

are working with as clearly as possible.1 (Midway et al).This article's presentation of an infographic detailing the phases and guiding

principles of data visualisation was quite effective. Two principles are indicated in yellow at the beginning of the design phase of a figure. The six green principles serve as an overview of the judgements and general considerations that must be taken when creating a figure. The two guidelines in blue are generally thought of as the last steps once a figure has been generated. Even though the basic flow of the principles is from bottom to top, it is likely that distinct principles will need to be taken into account in a different sequence when constructing specific figures (Midway et al, 2020).

## Collecting Data:

Collecting a lot of data is the first and most crucial stage in data visualisation. We can only use data visualisation techniques to the gathered data and gain some useful insights from it once we have a significant amount of data.

## Clean Your Data:

Prior to constructing a visualization, data cleansing must be carried out. A huge dataset may contain a number of data points with unsuitable, erroneous, or fake values that could result in the addition of anomalous graphics.

The output of a data cleaning procedure is often a dataset devoid of mistakes, anomalies, etc., which provides significantly greater accuracy when processing data. The dataset domain you're working with largely determines how you clean data.

## Choose Chart Type:

It is critical to understand your audience before selecting a visual chart or graph that would effectively express the content.

The type of chart you use is determined by the information you need to convey to your audience.

**Do you want to demonstrate how combining data columns can provide valuable insights?**

**Would you like to display some data patterns from the datasets?**

**Do you wish to demonstrate how data variables are compared to one another?**

**Do you wish to display the data variables' relationships?**

Choosing a few of them can assist you in determining which charts are best for you. This normally necessitates some trial and error with various charts before settling on the optimal one.

**Prepare Data:**

To prepare the data for visualization, determine the sort of graph, chart, or other visualisations you need to build, as well as the supporting library you will be integrating. After the chart is completed, the data may need to be transformed to meet the specifications.

Finding data columns that help make decisions, providing relevant insights into data, grouping data, establishing aggregate values for groups, combining variables to form new columns, and so on are all examples of data preparation chores.

**Visualize Data:**

Finally, you'll have the data you need to make visualisations. We now apply visualisation abilities to the prepared data and depict it in charts or graphs that provide significant insights. We'll present it to you after we're completed so you can see what we've discovered!.


**Data:**

The Global Health Observatory (GHO) data repository of the World Health Organisation (WHO) keeps track of various pertinent factors, including each country's health status. The datasets are made public so that they can be used by anybody to analyse health data. The dataset for life expectancy, health factors for 193 countries, and associated economic figures were all sourced from the same WHO data repository website. From each category of health-related issues,

only the most representative critical factors were chosen. Human mortality rates have improved significantly over the past 15 years, especially in emerging nations, when compared to the previous 30 years, it has been highlighted. Data for 193 countries were used from the years 2000 to 2015 to conduct additional analysis for this project. The various data sets were combined to form a single dataset. Some values were missing when the data were first visually inspected. Because the datasets were provided by the WHO, we did not discover any obvious errors. To manage missing data, the Missmap command in R was employed. The population, Hepatitis B, and GDP had the most missing data, the results showed. The missing information was contributed by lesser-known countries like Vanuatu, Tonga, Togo, Cabo Verde, etc. It was decided to exclude these countries from the final model dataset since it was difficult to find all the data for these countries. The resulting merged file, which includes 22 columns and 2938 rows in the final dataset, contains 20 predicting variables. Then, all anticipated variables were divided into four main groups: immunisations, economic factors, societal issues, and mortality-related factors.

The publicly available dataset provides data for 193 countries spanning from year 2000 to year 2015 and is structured in **2938 rows** (data points) which are characterized into a total of **22 columns** (features). The features can be categorized into two groups:

- Health factors which are originally provided by the Global Health Observatory (GHO) data repository under the Wor Health Organization (WHO)
- Economic factors which have been collected by the United Nation (UN) website.

The link to the dataset is as below:

[Life Expectancy (WHO) | Kaggle](#)

This data contains the life expectancies of human population including both the infants as well as the adults living in the different regions of the world.

This data will provide the various insights regarding the current trends of the rise or fall in the life expectancy and the impacts on the GDP due to the mortality in the different part of the world.

```
Source  Visual                                                                              ☰ Outline

 1▾ ---
 2  title: "Data Visualisation"
 3  output: html_document
 4  date: "2023-04-14"
 5▴ ---
 6▾ ```{r}
 7  library(tidyverse)
 8  library(ggplot2)
 9  library(ggcorrplot)
10  library(corrplot)
11  library(leaps)
12  library(car)
13  library(Metrics)
14  library(reshape2)
15  library(ggpubr)
16  library(moments)
17▴ ```
18  Loading of the Life Expectancy Data Which i downloaded from a Website
19
20▾ ```{r}
21  data <- read.csv("./Life Expectancy Data.csv")
22  head(data)
23  sprintf("Dataset size: [%s]", toString(dim(data)))
24
25▴ ```
26
27▾ ```{r}
28  summary(data)
29▴ ```
30
31  Clean and filter data.
32  Missing data:There are no missing values in the data that i have loaded.
33
34▾ ```{r}
35  missing.rows = dim(data)[1] - dim(na.omit(data))[1]
36  sprintf("Dataset size: [%s]", toString(dim(data)))
37  sprintf("Missing rows: %s (%s%%)", missing.rows, round((missing.rows*100)/dim(data)[1], 2))
38
39  missings_df <- data.frame(type=c("missing", "non-missing"), count = c(missing_rows, dim(na.omit(data))[1]))
```

```
33
34 ▾ ```{r}
35   missing.rows = dim(data)[1] -  dim(na.omit(data))[1]
36   sprintf("Dataset size: [%s]", toString(dim(data)))
37   sprintf("Missing rows: %s (%s%%)", missing.rows, round((missing.rows*100)/dim(data)[1], 2))
38
39   missings_df <- data.frame(type=c("missing", "non-missing") ,count = c(missing.rows,  dim(na.omit(data))[1]))
40   missing_counts <- data.frame(feature = factor(names(data)),
41                       counts=sapply(data, function(x) sum(is.na(x))))
42
43 ▴ ```
44
45 ▾ ```{r}
46   sum(is.na(data))
47 ▴ ```
```

```
[1] 2563
```

```
48
49   Duplicate Values:The are also no duplicate valurs present in the datatset
50 ▾ ```{r}
51   sum(duplicated(data))
52 ▴ ```
```

```
[1] 0
```

```
53 ▾ ```{r}
54   laptops<- unique(data) #removing duplicated rows
55   sum(duplicated(data))
56 ▴ ```
```

```
[1] 0
```

```
57
58
59   My final summary.
```

```
57
58
59   My final summary.
60 ▾ ```{r}
61   summary(data)
62 ▴ ```
63
64
65
66
67   Finally, I am writing my file as a csv, in order to be able to have it also as an excel file.
68 ▾ ```{r}
69   write.csv(laptops, file = 'Life_Expectancy_Modified_Data.csv')
70 ▴ ```
71
72
```

## Data Dictionary:

The data dictionary shows the completed dataset that I am going to use, along with the variables that I will need to accomplish to answer my questions shown in bold:

| Variable | Type | Description |
|---|---|---|
| Country | factor | Country name |
| Year | numeric | Year of the data |
| Status | factor | Country status of developed or developing |
| Life_Expectancy | numeric | Life expectancy in age |
| Adult_Mortality | numeric | Adult Mortality Rates of both sexes (probability of dying between 15 and 60 years per 1000 population) |
| infant.deaths | numeric | Number of Infant Deaths per 1000 population |
| Alcohol | numeric | Alcohol, recorded per capita (15+) consumption (in litres of pure alcohol) |
| percentage.expenditure | numeric | Expenditure on health as a percentage of Gross Domestic Product per capita(%) |
| Hepatitis.B | numeric | Hepatitis B (HepB) immunization coverage among 1-year-olds (%) |
| Measles | numeric | number of reported cases per 1000 population |
| BMI | numeric | Average Body Mass Index of entire population |
| under.five.deaths | numeric | Number of under-five deaths per 1000 population |
| Polio | numeric | Polio (Pol3) immunization coverage among 1-year-olds (%) |
| Total.expenditure | numeric | General government expenditure on health as a percentage of total government expenditure (%) |
| Diphtheria | numeric | Diphtheria tetanus toxoid and pertussis (DTP3) immunization coverage among 1-year-olds (%)room) |
| HIV.AIDS | numeric | Deaths per 1 000 live births HIV/AIDS (0-4 years) |
| GDP | numeric | Gross Domestic Product per capita (in USD) |
| Population | numeric | Population of the country |
| thinness..1.19.years | numeric | Prevalence of thinness among children and adolescents for Age 10 to 19 (% ) |
| thinness.5.9.years | numeric | Prevalence of thinness among children for Age 5 to 9(%) |
| Income.composition.of.resources | numeric | Human Development Index in terms of income composition of resources (index ranging from 0 to 1) |
| Schooling | numeric | Number of years of Schooling(years) |

## Persona And Questions:

The numerous visualisations are created in order to comprehend the diverse life expectancies in various parts of the world.

We will create a dashboard which will depict all the various dependencies and the analysis.

Questions

Q1: What is the average life expectancy across various regions of the world?

Q2: How does the GDP get affected by the Life Expectancy?

Q3: What is the trend of the life expectancy over the past years?

Q4: What is the maximum life Expectancy for each country?

Q5: What are the total Deaths of infants according to the region(Also find the diseases which caused death)?

**Requirements:**
Specifying the main requirements in terms of the relationships that must be visualised to answer each question, followed by brief design ideas for representations and interactions.


**Design:**

It can be useful to distinguish between "top-down" and "bottom-up" visual attention when attempting to determine which aspects of a data visualisation attract and which do not. What is referred described as "top-down" visual focus is determined by viewer expectations and objectives. Contrarily, physical aspects of the image, such colour and contrast, are what command bottom-up visual attention(2017 ,Matzen et al). In addition, I also employ a layout technique known as the z-pattern that imitates the path the reader's eye follows as it moves from left to right and top to bottom.

For the researchers of WHO and other organisations to know  the reasons for the decrease in the life expectancy of human and the reasons of the reduction in life expectancy, Data visualisation plays an active role and it generates accurate results that can be used for the purpose of evaluation.

Today's business environment benefits from data visualisation since it helps senior management analyse huge and complicated data sets. This procedure is made easier by interactive data visualisation (IDV), which gives consumers a straightforward interface to browse, select, and visualise data. Data analytics typically includes the usage of IDV (Janvrin et al., 2014).

**Implementation:**

After writing the csv file with RStudio, I open the excel file and I delete the first column that has been created automatically which is the count of rows.

After saving it, I import it into Tableau.



To answer the above mentioned questions,I need to create the various visualisations required.

## Q1: What is the average life expectancy across various regions of the world?
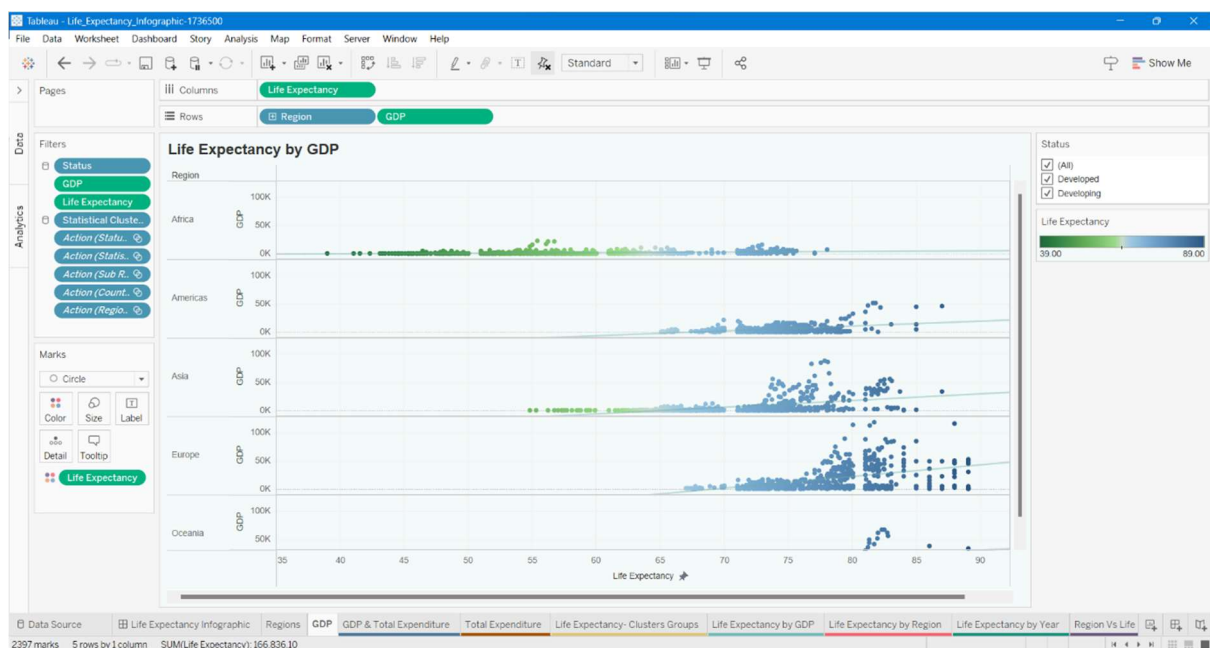


The visualisation is plotted between the fields such as the regions and the average life expectancy of human population.

The blue regions show the highest life expectancy and the areas that are towards the red show the lowest life expectancy.

The highest average life expectancy was noted in Europe and the lowest in Africa.

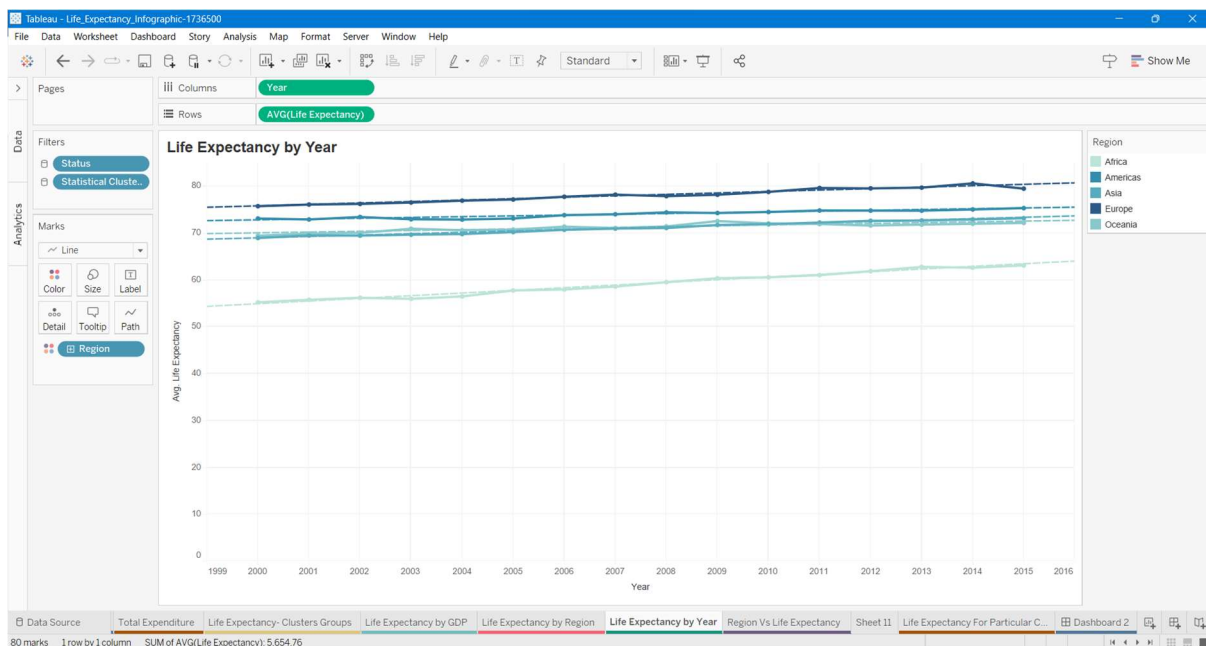## Q2: How does the GDP get affected by the Life Expectancy?



This visualisation indicated the impact on the GDP due to the life expectancy.

The visual shows the various regions or continents and the total GDP how it is related to the life expectancy.

Here I have used the scatter plot for plotting my data and added the filters according to the need of visualisation.

The GDP can be filtered on the basis of developing and developed countries in a particular region. Further more the more darker the bubbles the more life expectancy observed.

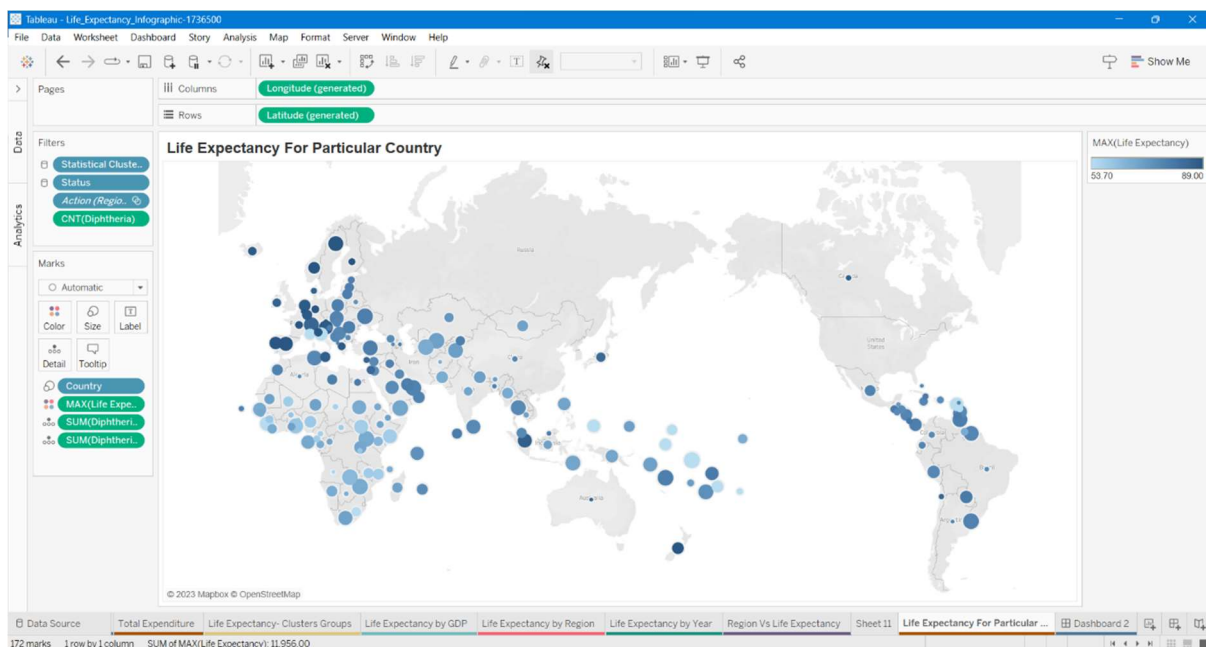## Q3: What is the trend of the life expectancy over the past years?



Here we have plotted the graph between the average life expectancy to the Years.

This graph shows the trend of the life expectancy as to whether it is increasing or decreasing over the past and the upcoming years.

Tableau does some statistical calculations and visualises the trend line as to whether the life expectancy will increase of decrease for the upcoming years.
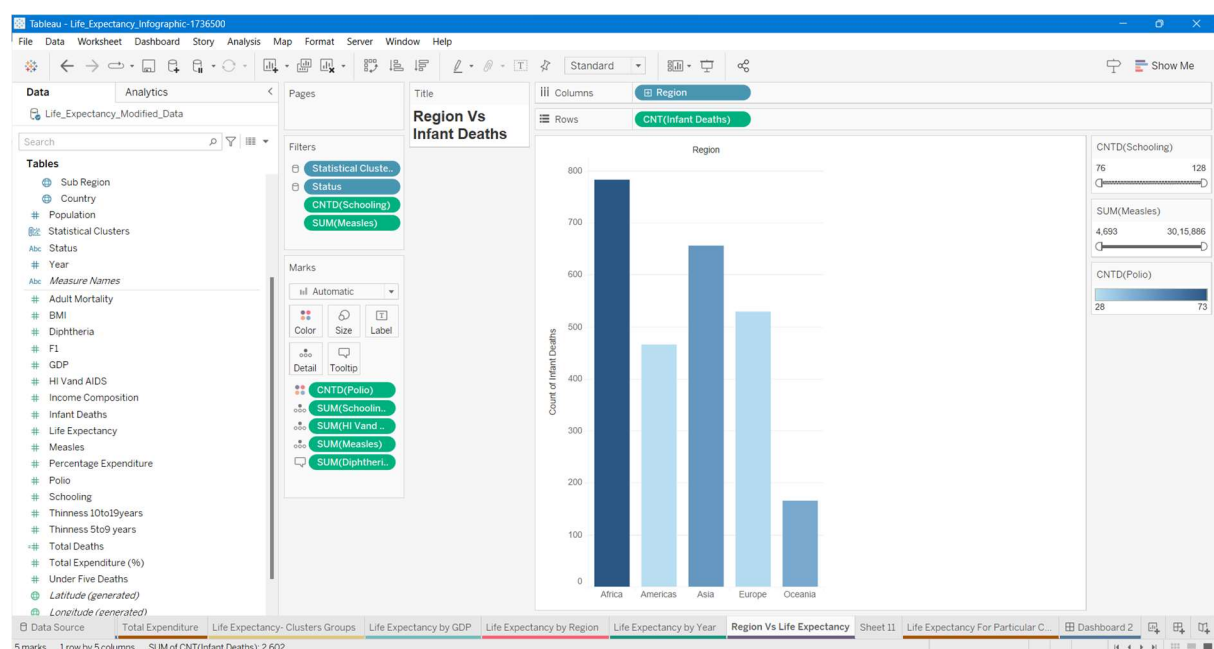
## Q4: What is the maximum life Expectancy for each country?

For this visualisation, I have used a map for illustrating the maximum life expectancy across different countries.

This visualisation show the maximum life expectancy for each individual country. The more darker the bubble the more life expectancy observed for that particular country.

**Q5: What are the total Deaths of infants according to the region(Also find the diseases which caused death)?**
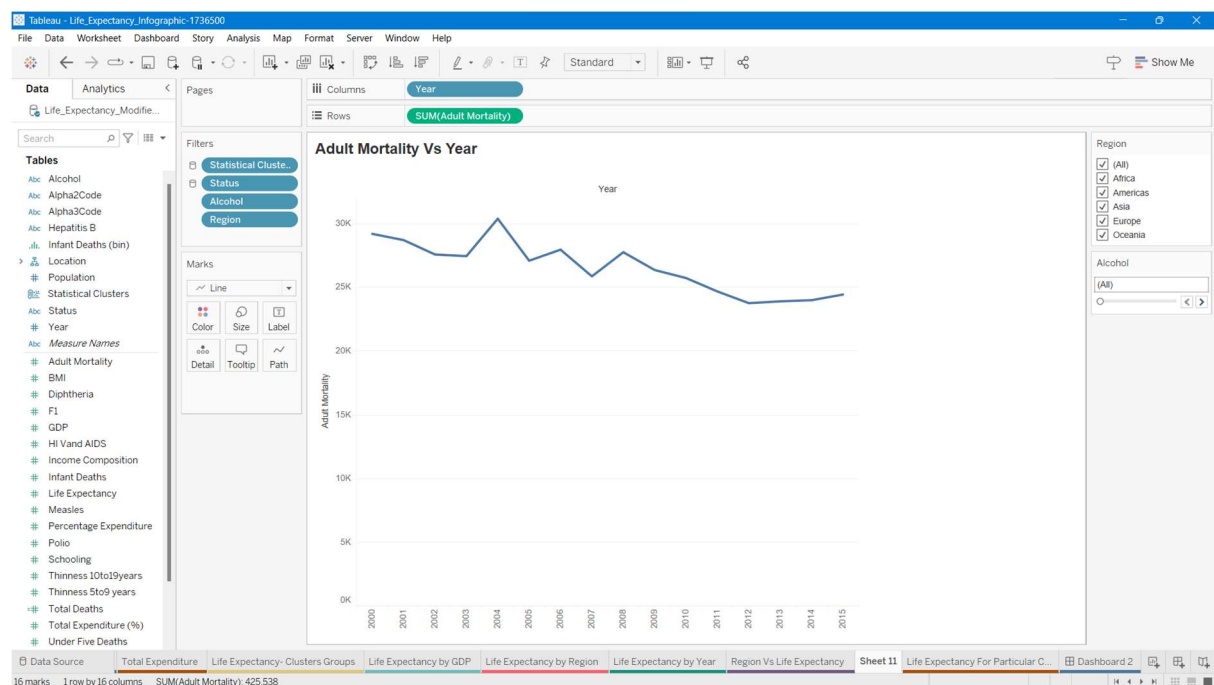


Here we have plotted the bar graph for the rate of mortality of infants according to the regions of the world and also filtered the data according to the children who were schooling.

**Mortality (Adults) vs Year:**

The line chart illustrates the sum of mortality over the years from 2000 to 2015 and filtered the data according the to the mortality due to alcohol consumption and got the analysis.

The graph was a descending one as the adult mortality rates decreased as the years passed and there is a negative trend towards the rise in the adult mortality

Also added the filter according to the region to get better insights on to which region has the highest mortality rate due to alcohol consumption.



**Dashboards:**

A dashboard is a way of displaying various types of visual data in one place. Usually, a dashboard is intended to convey different, but related information in an easy-to-digest form. And oftentimes, this includes things like key performance indicators (KPI)s or other important business metrics that stakeholders need to see and understand at a glance.
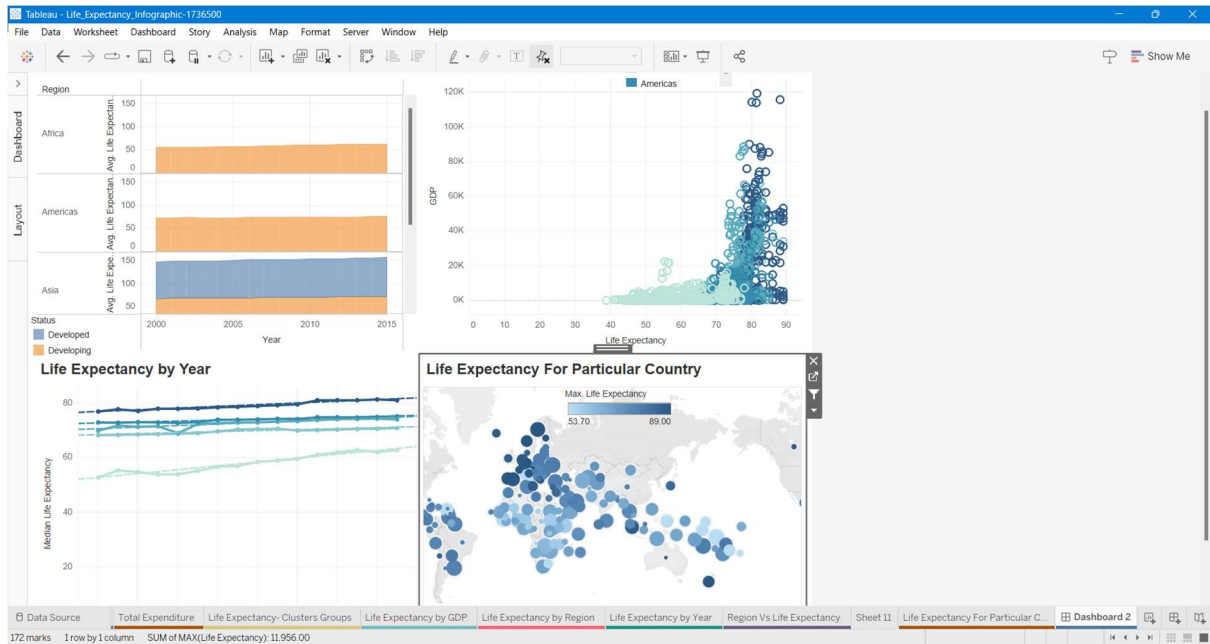
**Data dashboards**:

They are a summary of different, but related data sets, presented in a way that makes the related information easier to understand. Dashboards are a type of data visualization, and often use common visualization tools such as graphs, charts, and tables.

An interactive dashboard was made for the purpose of getting valuable insights from all the data and filtering them according to the other graphs plotted.
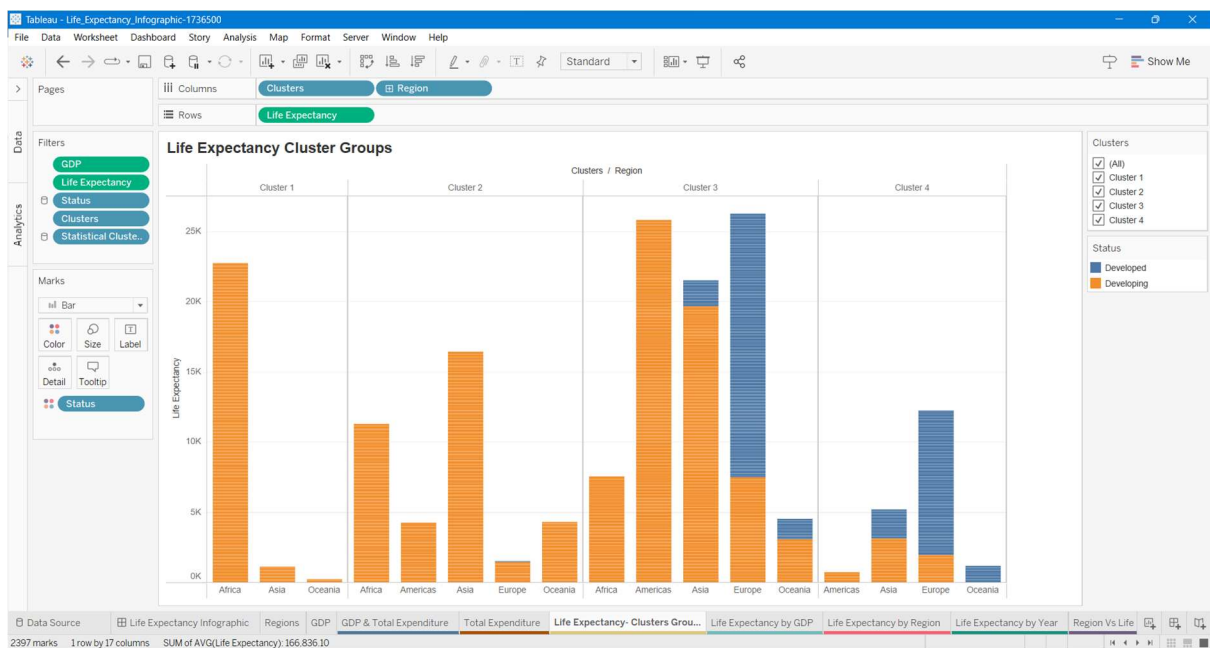
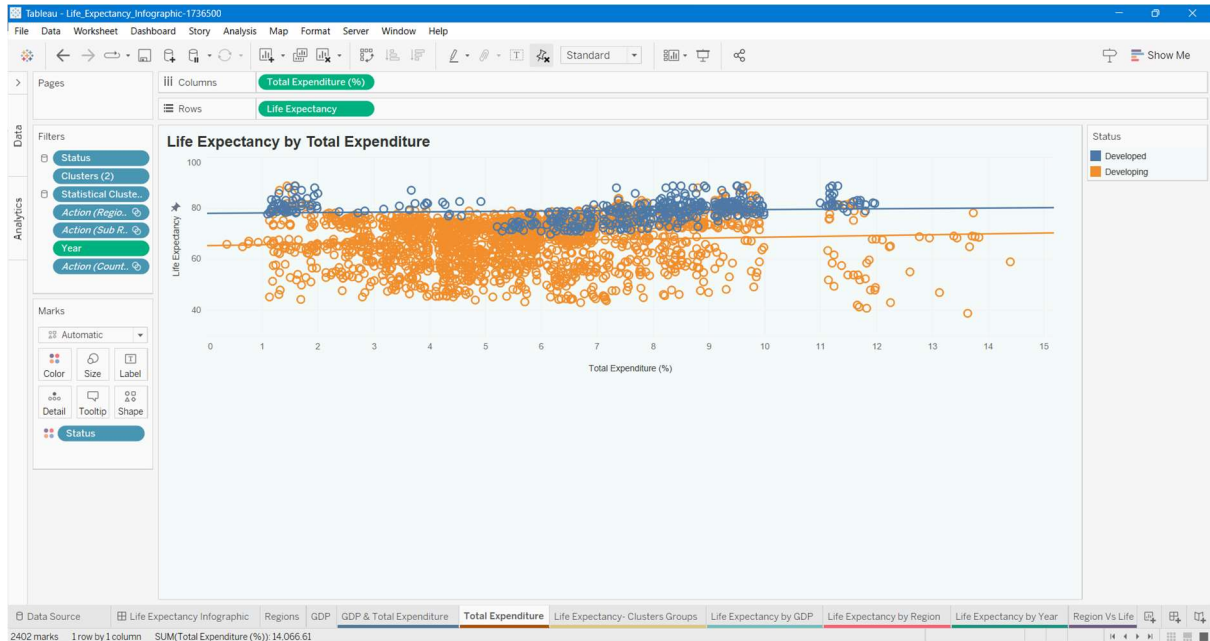## Second Interactive Dashboard:



The other plotted graphs are as below:
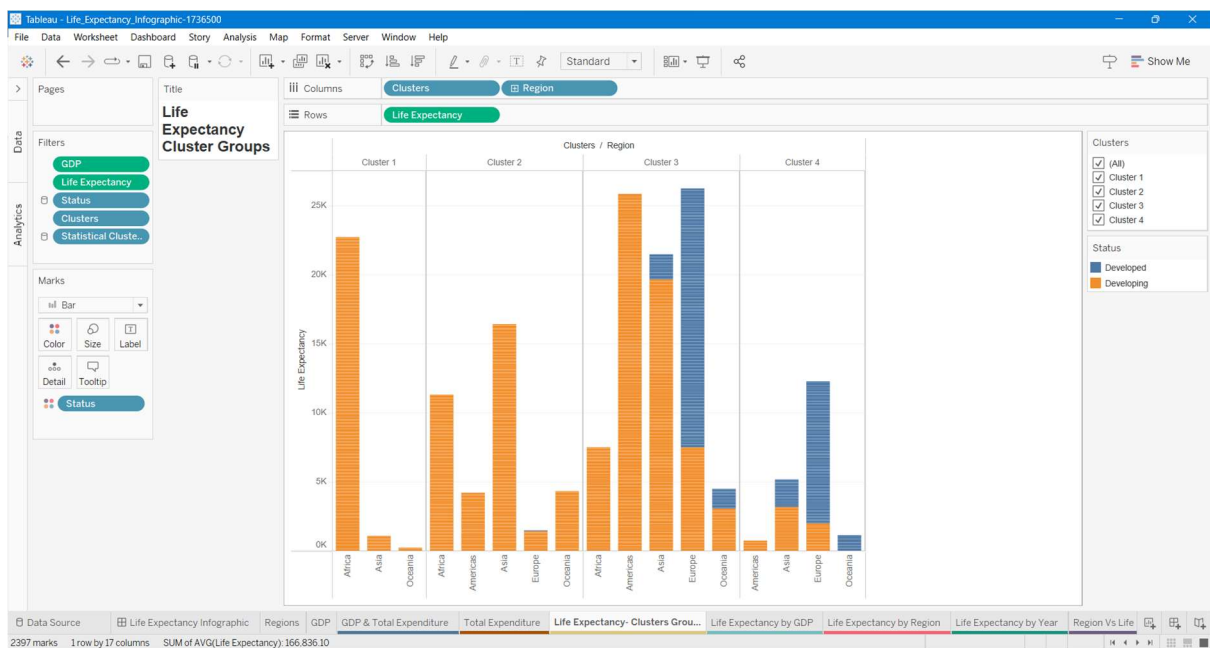
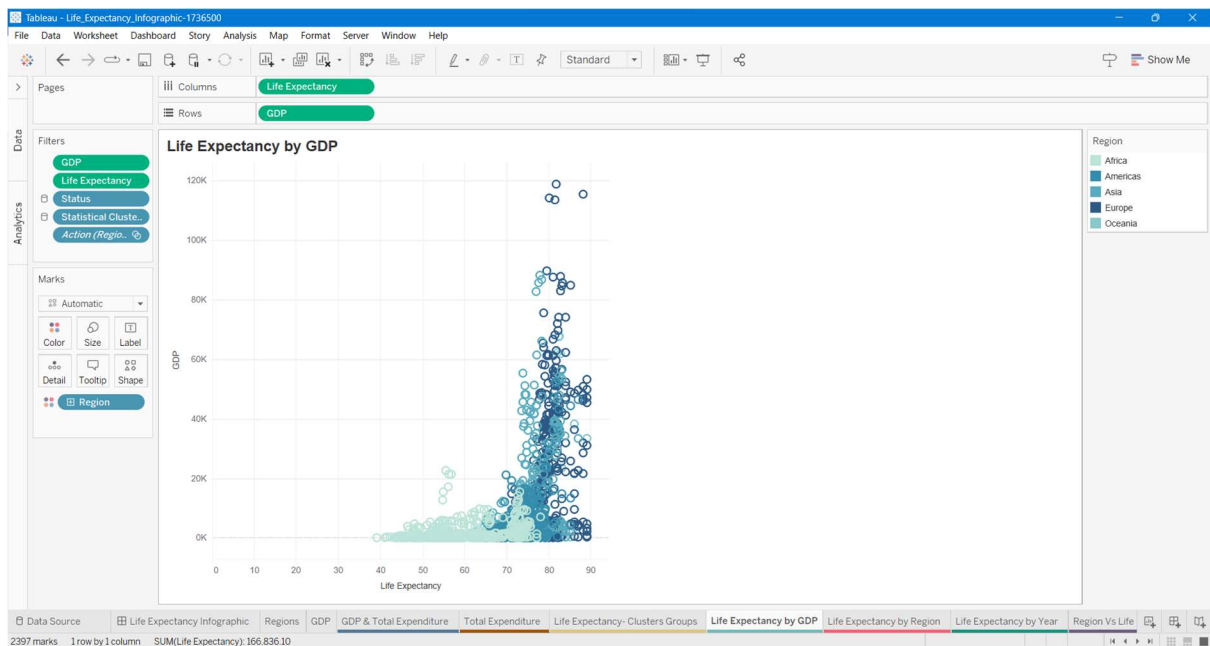## Life Expectancy Cluster Groups:

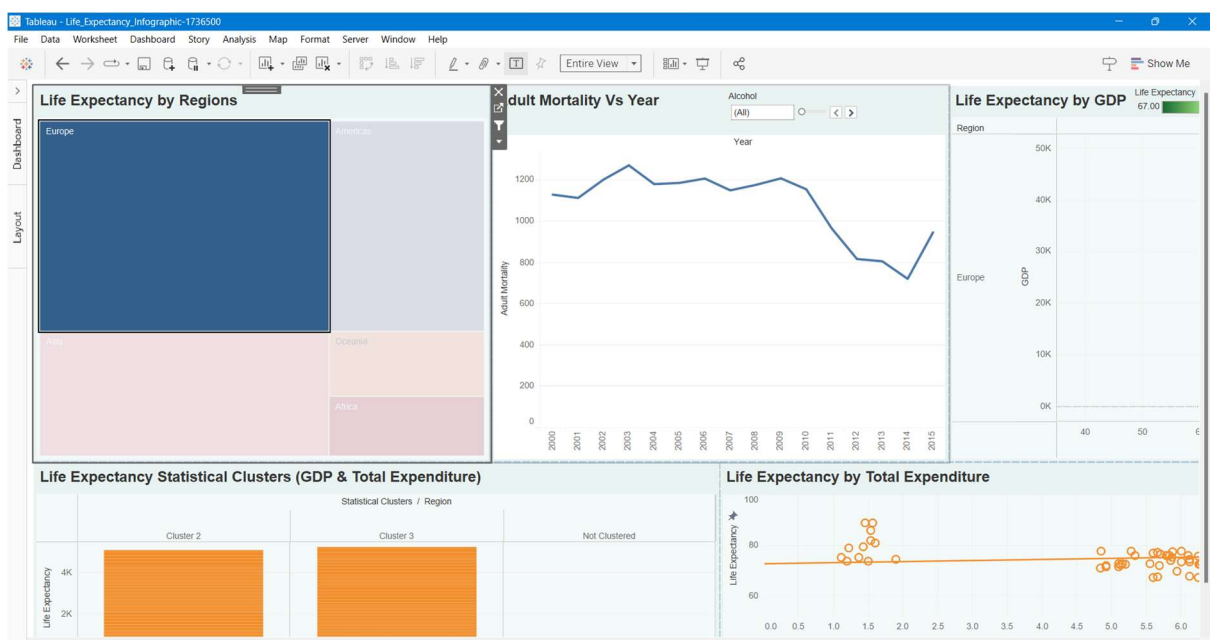# Life Expectancy by Total Expenditure:



# Life Expectancy Cluster Groups

# Life Expectancy by GDP



## WalkThrough:



The data gets filtered according to the selection of a particular region from any chart and we can observe a change in all the charts having a change according the filter done.

These dashboards will help the organisation to solve many problems that are prevailing according to the regions or country wise.

For Example:

A officer clicks on the Europe region, then he would be able to view all the life expectancies throughout Europe and would be able to know on what to do to solve the issue.

It can be observed that the adult mortality rate has quite decrease from the previous years while also taking a sight at the GDP too.

**Reflective Discussion:**

To begin, I learned from this that Exploratory Data Analysis is one of the most critical skills that a Data Scientist must possess to comprehend a dataset. Prior to exploring this module, EDA was a time-consuming procedure that necessitated considerable code in either R or Python.

However, after learning how to utilise Tableau, I realised that this approach could be greatly shortened, as Tableau is a very user-friendly framework that requires little acclimating to grasp the fundamentals. Apart from that, Tableau can make use of code to provide the user with more options, that can be accessed only with using mouse and not writing code, if the user prefers. Something that in Power BI wasn't so easy. Initially, I found it more complicated and, in my opinion, more limited than Tableau, but I realised that, aside from the possibility of a higher grade, the most important thing that I believe will benefit me in general is that, after completing this assignment, I will be able to include not only Tableau but also Power BI on my CV, as I have now had a taste of it. This is critical, because many businesses, as far as I am aware, want one of these tools.

Finally, the most significant limitation that I can think of is that these tools aren't free, and in the absence of licence, only very basic graphs are offered, in comparison to R and Python, for example.

**References:**

1. Azzam, T., Evergreen, S., Germuth, A.A. and Kistler, S.J., 2013. Data visualization and evaluation. New Directions for Evaluation, 2013(139), pp.7-32.
2. Midway, S.R., 2020. Principles of effective data visualization. *Patterns*, *1*(9), p.100141.
3. Matzen, L.E., Haass, M.J., Divis, K.M. and Stites, M.C., 2017, July. Patterns of attention: How data visualizations are read. In *International Conference on Augmented Cognition* (pp. 176-191). Springer, Cham.
4. Marr, B., 2012. *Key Performance Indicators (KPI): The 75 measures every manager needs to know*. Pearson UK.
5. Janvrin, D.J., Raschke, R.L. and Dilla, W.N., 2014. Making sense of complex data using interactive data visualization. *Journal of Accounting Education*, *32*(4), pp.31-48.