**FLIP ROBO**

# CAUSE OF DEATH

**Prepared by: Maahi Gurnani, Data Science Intern at Flip Robo Technologies**

# ACKNOWLEDGEMENT

# Introduction

- ## **Problem Statement:**

A straightforward way to assess the health status of a population is to focus on mortality – or concepts like child mortality or life expectancy, which are based on mortality estimates. A focus on mortality, however, does not take into account that the burden of diseases is not only that they kill people, but that they cause suffering to people who live with them. Assessing health outcomes by both mortality and morbidity (the prevalent diseases) provides a more encompassing view on health outcomes. This is the topic of this entry. The sum of mortality and morbidity is referred to as the 'burden of disease' and can be measured by a metric called 'Disability Adjusted Life Years' (DALYs). DALYs are measuring lost health and are a standardized metric that allow for direct comparisons of disease burdens of different diseases across countries, between different populations, and over time. Conceptually, one DALY is the equivalent of losing one year in good health because of either premature death or disease or disability. One DALY represents one lost year of healthy life. The first 'Global Burden of Disease' (GBD) was GBD 1990 and the DALY metric was prominently featured in the World Bank's 1993 World Development Report. Today it is published by both the researchers at the Institute of Health Metrics and Evaluation (IHME) and the 'Disease Burden Unit' at the World Health Organization (WHO), which was created in 1998. The IHME continues the work that was started in the early 1990s and publishes the Global Burden of Disease study.

- ## **Motivation for the Problem Undertaken:**

Measuring how many people die each year and why they have died is one of the most informative ways of assessing the effectiveness of a country's health system. Mortality data allow health authorities to evaluate how they prioritize public health programs.

Aggregation was done over estimates of "Total Deaths", "Deaths by Cause" and by "Country/Territory" to estimate regional and global cause specific mortality rates.

# Analytical Problem Framing

- **Analytical Modelling of the Problem:**

  1. Checking the year wise count of total deaths from 1990-2019
     -Displaying Total Deaths of all Countries in Descending order

  2. Analysis with respect to either
     a) Each Disease that has affected which Top 10 Countries 'or'
     b) Top 10 Diseases of each Country

# Hardware and Software Requirements and Tools Used

- **Languages Used:** Python

- **Platform Used:** Jupyter Notebook

- **Libraries and Metrics used:**

Following are the libraries used to carry our analysis

```
In [1]: #importing the necessary libraries
        import numpy as np
        import pandas as pd
        import seaborn as sns
        import matplotlib.pyplot as plt
        import warnings
```

# Data Sources and their formats

In this Dataset, we have Historical Data of different cause of deaths for all ages around the World. The key features of this Dataset are: Meningitis, Alzheimer's Disease and Other Dementias, Parkinson's Disease, Nutritional Deficiencies, Malaria, Drowning, Interpersonal Violence, Maternal Disorders, HIV/AIDS, Drug Use Disorders, Tuberculosis, Cardiovascular Diseases, Lower Respiratory Infections, Neonatal Disorders, Alcohol Use Disorders, Self-harm, Exposure to Forces of Nature, Diarrheal Diseases, Environmental Heat and Cold Exposure, Neoplasms, Conflict and Terrorism, Diabetes Mellitus, Chronic Kidney Disease, Poisonings, Protein-Energy Malnutrition, Road Injuries, Chronic Respiratory Diseases, Cirrhosis and Other Chronic Liver Diseases, Digestive Diseases, Fire, Heat, and Hot Substances, Acute Hepatitis.

Including the snapshot of the data set provided and loaded

```
In [2]: df_report=pd.read_csv("cause of death.csv")
        df_report
```

Out[2]:

| | Country/Territory | Code | Year | Meningitis | Alzheimer's Disease and Other Dementias | Parkinson's Disease | Nutritional Deficiencies | Malaria | Drowning | Interpersonal Violence | ... | Diabetes Mellitus | Chronic Kidney Disease | Poisonings |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | AFG | 1990 | 2159 | 1116 | 371 | 2087 | 93 | 1370 | 1538 | ... | 2108 | 3709 | 338 |
| 1 | Afghanistan | AFG | 1991 | 2218 | 1136 | 374 | 2153 | 189 | 1391 | 2001 | ... | 2120 | 3724 | 351 |
| 2 | Afghanistan | AFG | 1992 | 2475 | 1162 | 378 | 2441 | 239 | 1514 | 2299 | ... | 2153 | 3776 | 386 |
| 3 | Afghanistan | AFG | 1993 | 2812 | 1187 | 384 | 2837 | 108 | 1687 | 2589 | ... | 2195 | 3862 | 425 |
| 4 | Afghanistan | AFG | 1994 | 3027 | 1211 | 391 | 3081 | 211 | 1809 | 2849 | ... | 2231 | 3932 | 451 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 6115 | Zimbabwe | ZWE | 2015 | 1439 | 754 | 215 | 3019 | 2518 | 770 | 1302 | ... | 3176 | 2108 | 381 |
| 6116 | Zimbabwe | ZWE | 2016 | 1457 | 767 | 219 | 3056 | 2050 | 801 | 1342 | ... | 3259 | 2160 | 393 |
| 6117 | Zimbabwe | ZWE | 2017 | 1460 | 781 | 223 | 3090 | 2116 | 819 | 1363 | ... | 3313 | 2196 | 398 |

# Data Inputs- Logic- Output Relationships

There are total 34 columns, of which 31 are the various causes of disease and 6120 rows containing 30 years data(1990-2019) of 204 Countries. Below are the details of all the input variables and their entries.

- 01. Country/Territory - Name of the Country/Territory
- 02. Code - Country/Territory Code
- 03. Year - Year of the Incident
- 04. Meningitis - No. of People died from Meningitis
- 05. Alzheimer's Disease and Other Dementias - No. of People died from Alzheimer's Disease and Other Dementias
- 06. Parkinson's Disease - No. of People died from Parkinson's Disease
- 07. Nutritional Deficiencies - No. of People died from Nutritional Deficiencies
- 08. Malaria - No. of People died from Malaria
- 09. Drowning - No. of People died from Drowning
- 10. Interpersonal Violence - No. of People died from Interpersonal Violence
- 11. Maternal Disorders - No. of People died from Maternal Disorders
- 12. Drug Use Disorders - No. of People died from Drug Use Disorders
- 13. Tuberculosis - No. of People died from Tuberculosis
- 14. Cardiovascular Diseases - No. of People died from Cardiovascular Diseases
- 15. Lower Respiratory Infections - No. of People died from Lower Respiratory Infections
- 16. Neonatal Disorders - No. of People died from Neonatal Disorders
- 17. Alcohol Use Disorders - No. of People died from Alcohol Use Disorders
- 18. Self-harm - No. of People died from Self-harm
- 19. Exposure to Forces of Nature - No. of People died from Exposure to Forces of Nature
- 20. Diarrheal Diseases - No. of People died from Diarrheal Diseases
- 21. Environmental Heat and Cold Exposure - No. of People died from Environmental Heat and Cold Exposure
- 22. Neoplasms - No. of People died from Neoplasms
- 23. Conflict and Terrorism - No. of People died from Conflict and Terrorism
- 24. Diabetes Mellitus - No. of People died from Diabetes Mellitus
- 25. Chronic Kidney Disease - No. of People died from Chronic Kidney Disease
- 26. Poisonings - No. of People died from Poisoning
- 27. Protein-Energy Malnutrition - No. of People died from Protein-Energy Malnutrition
- 28. Chronic Respiratory Diseases - No. of People died from Chronic Respiratory Diseases
- 29. Cirrhosis and Other Chronic Liver Diseases - No. of People died from Cirrhosis and Other Chronic Liver Diseases
- 30. Digestive Diseases - No. of People died from Digestive Diseases
- 31. Fire, Heat, and Hot Substances - No. of People died from Fire or Heat or any Hot Substances
- 32. Acute Hepatitis - No. of People died from Acute Hepatitis

# Data Pre-processing

1. **Checking for null values:** We should deal with the problem of missing values

**df_report.isnull().sum()**

```
In [6]: #Checking for null values
        df_report.isnull().sum()

Out[6]: Country/Territory                              0
        Code                                           0
        Year                                           0
        Meningitis                                     0
        Alzheimer's Disease and Other Dementias        0
        Parkinson's Disease                            0
        Nutritional Deficiencies                       0
        Malaria                                        0
        Drowning                                       0
        Interpersonal Violence                         0
        Maternal Disorders                             0
        HIV/AIDS                                       0
        Drug Use Disorders                             0
        Tuberculosis                                   0
        Cardiovascular Diseases                        0
        Lower Respiratory Infections                   0
        Neonatal Disorders                             0
        Alcohol Use Disorders                          0
        Self-harm                                      0
```

Obs- There is no missing data in the dataset, so no need to treat it

2. **Checking for Country Count**

   **df_report['Country/Territory'].nunique()**

```
In [7]: #Checking the number of unique values in Country
        df_report['Country/Territory'].nunique()

Out[7]: 204
```

Obs- There is total 204 countries mortality data present in the dataset

3. **Checking for total number of Years Count**

   **df_report['Year'].nunique()**

```
In [8]: #Checking the number of unique values in Year
        df_report['Year'].nunique()

Out[8]: 30
```

Obs- We can see we have mortality-data of 30 years of each country ranging from 1990-2019

# Data Analysis and Visualization

For getting the insights of relationship between the various features we considered  doing certain analysis(We carried on our analysis with respect to either a)Each Disease that has affected which Top 10 Countries 'or' b) Top 10 Diseases of each Country throughout so as to get detailed insights of which country has been affected the most and with which disease.) and   visualization to discover any hidden patterns.

## 1.  Checking the year wise count of total deaths from 1990-2019

```
In [10]: df_report["total_death"] = df_report.iloc[:,3:].sum(axis=1)  # Since the cause of death is starting from 3rd column
         death_by_year = df_report.groupby("Year").sum()["total_death"]
         death_by_year
```

```
Out[10]: Year
         1990    43518516
         1991    44059729
         1992    44459130
         1993    45185713
         1994    46182613
         1995    46177018
         1996    46320827
         1997    46672370
         1998    47066088
         1999    47652090
         2000    48050317
         2001    48385692
         2002    48897031
         2003    49123952
         2004    49330171
         2005    49591909
         2006    49424521
         2007    49495216
         2008    50115740
         2009    49900666
         2010    50422775
         2011    50413303
         2012    50597654
         2013    50931550
         2014    51268375
         2015    51856393
         2016    52337435
         2017    52789758
         2018    53545244
         2019    54362920
         Name: total_death, dtype: int64
```

## 1.1 Plotting a line plot to check total deaths of all years(1990-2019)

```
In [11]: fig = plt.figure(figsize=(9,7))
         sns.lineplot(data=death_by_year,marker="o")
         plt.ylabel("No of Deaths")

Out[11]: Text(0, 0.5, 'No of Deaths')
```

Obs- We can see the number of deaths have increased over the years.

- **Further Analysis:**

We carried on our analysis with respect to either :

      a) Each Disease (Cause_of_death) that has affected which Top 10 Countries 'or'

      b) Top 10 Diseases of each Country

throughout so as to get detailed insights of which country has been affected the most and with which disease.

- **Steps Followed:**
    a) We want to take a look at the total global deaths and which countries lead in total deaths.
    b) Thus, we are aggregating across the DataFrame and would therefore drop the 'Year' column.
    c) And the Total Global Deaths from 1990-2019 will be displayed at end

```
In [12]: # drop the year column so we can aggregate the deaths from each cause without counting the year
         df2 = df_report.drop(labels='Year', axis=1).groupby(by='Country/Territory').sum().reset_index()
```

```
In [13]: df3=df2.drop(columns=['total_death'])
         df3
         # Dropped the self_created 'total-death' column and appending Total Deaths columns to get the sum of deaths of all diseases
```

Out[13]:

| | Country/Territory | Meningitis | Alzheimer's Disease and Other Dementias | Parkinson's Disease | Nutritional Deficiencies | Malaria | Drowning | Interpersonal Violence | Maternal Disorders | HIV/AIDS | ... | Diabetes Mellitus | Chronic Kidney Disease | Poisonin |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 78666 | 41998 | 13397 | 71453 | 13924 | 56536 | 108228 | 129621 | 4282 | ... | 93207 | 134676 | 14£ |
| 1 | Albania | 1323 | 16549 | 4491 | 569 | 0 | 2397 | 5242 | 246 | 57 | ... | 4055 | 7636 | £ |
| 2 | Algeria | 15685 | 86914 | 22943 | 7138 | 70 | 24273 | 16702 | 29475 | 6101 | ... | 89035 | 154666 | 12: |
| 3 | American Samoa | 30 | 143 | 69 | 60 | 0 | 120 | 101 | 30 | 15 | ... | 970 | 512 | |
| 4 | Andorra | 0 | 614 | 137 | 0 | 0 | 0 | 15 | 0 | 85 | ... | 198 | 292 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 199 | Venezuela | 11615 | 108735 | 18573 | 22554 | 3726 | 20273 | 266071 | 12739 | 46090 | ... | 175790 | 161667 | 2€ |
| 200 | Vietnam | 38559 | 369363 | 83322 | 48613 | 17311 | 214356 | 47981 | 13167 | 148838 | ... | 544222 | 396874 | 34€ |
| 201 | Yemen | 21095 | 31045 | 7188 | 68939 | 143463 | 27994 | 17918 | 53611 | 6276 | ... | 30812 | 52119 | 125 |
| 202 | Zambia | 98886 | 13473 | 4054 | 95913 | 205529 | 12809 | 30065 | 28395 | 1175563 | ... | 54098 | 41751 | 9( |
| 203 | Zimbabwe | 41238 | 20017 | 5764 | 66723 | 118728 | 18169 | 32741 | 29802 | 1836042 | ... | 71175 | 49952 | 9' |

204 rows × 32 columns

```
In [14]: df_report_mod = df3
         df_report_mod['Total Deaths'] = df3.sum(axis=1)
         df_report_mod.head()
```

Out[14]:

| | Drowning | Interpersonal Violence | Maternal Disorders | HIV/AIDS | ... | Chronic Kidney Disease | Poisonings | Protein-Energy Malnutrition | Road Injuries | Chronic Respiratory Diseases | Cirrhosis and Other Chronic Liver Diseases | Digestive Diseases | FireHeat and Hot Substances | Acute Hepatitis | Total Deaths |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 56536 | 108228 | 129621 | 4282 | ... | 134676 | 14530 | 70163 | 208331 | 209857 | 98419 | 186959 | 13559 | 98108 | 5982631 |
| | 2397 | 5242 | 246 | 57 | ... | 7636 | 500 | 526 | 8522 | 22632 | 8717 | 14907 | 636 | 44 | 523545 |
| | 24273 | 16702 | 29475 | 6101 | ... | 154666 | 12337 | 6407 | 369395 | 168453 | 91927 | 146527 | 27628 | 10492 | 4601205 |
| | 120 | 101 | 30 | 15 | ... | 512 | 0 | 60 | 164 | 612 | 181 | 341 | 0 | 0 | 8619 |
| | 0 | 15 | 0 | 85 | ... | 292 | 0 | 0 | 259 | 838 | 283 | 560 | 0 | 30 | 12532 |

## 2. Displaying Total Deaths of all Countries in Descending order

```
In [15]: country_wise = df_report_mod[['Country/Territory', 'Total Deaths']].sort_values(by='Total Deaths', ascending=False).reset_in
         country_wise
```

Out[15]:

| | Country/Territory | Total Deaths |
|---|---|---|
| 0 | China | 265408106 |
| 1 | India | 238158165 |
| 2 | United States | 71197802 |
| 3 | Russia | 59591155 |
| 4 | Indonesia | 44046941 |
| ... | ... | ... |
| 199 | Cook Islands | 3999 |
| 200 | Tuvalu | 2962 |
| 201 | Nauru | 2249 |
| 202 | Niue | 591 |
| 203 | Tokelau | 299 |

Obs- We can see the maximum mortality rate is of China while the least is of Tokelau, (However, population too plays important role, since more the population, more would be the mortality rate)

# 3. Displaying Top 10 countries with max no. of deaths

```
In [16]: top_10_country_deaths = country_wise[:10]
         top_10_country_deaths
```

Out[16]:

|   | Country/Territory | Total Deaths |
|---|---|---|
| 0 | China | 265408106 |
| 1 | India | 238158165 |
| 2 | United States | 71197802 |
| 3 | Russia | 59591155 |
| 4 | Indonesia | 44046941 |
| 5 | Nigeria | 43670014 |
| 6 | Pakistan | 38151878 |
| 7 | Brazil | 32674112 |
| 8 | Japan | 31922807 |
| 9 | Germany | 25559667 |

Obs- The top 10 countries with max mortality rates are 1)China, 2) India, 3)United States, 4)Russia 59591155, 5)Indonesia 44046941, 6)Nigeria, 7) Pakistan, 8) Brazil, 9) JapaN, 10) German

## 3.1    Plotting a bar plot to view top 10 Countries death rates

```
In [17]: # Plotting bar plot for top 10 mortality rate countries
         plt.figure(figsize=(8,8))
         plt.bar(top_10_country_deaths['Country/Territory'],top_10_country_deaths['Total Deaths'])
         plt.xlabel("Countries")
         plt.ylabel("Mortality rate")
         plt.title("Top 10 mortality rate countries")
         plt.xticks(rotation ='vertical')
         plt.show()
```

# 4. Displaying Every Cause of Death and top 10 Countries suffering from it

- We have taken every cause of death and have extracted the top 10 Countries (Sort function) suffering from that particular disease and plotted bar plots for the same

```
In [18]: for cause in df_report_mod.columns[1:]:    # Since the cause of death starts from 0 Index
             # Retrieves the top 10 countries/territories with the highest overall deaths spanning all years
             data = df_report_mod.set_index('Country/Territory')[cause].sort_values(ascending=False)[:10]

             # configurations for the bar graphs
             plt.figure(figsize=(10, 6))
             plt.bar(data = data, x = data.index, height = data.values, width=.6, color=['blue','red','green','magenta','yellow','cyan','b
             plt.xticks(rotation = 90)
             plt.xlabel('Countries', size= 13)
             plt.ylabel(cause +' Deaths')
             plt.title('Countries w/ the Most ' + cause +' Deaths')
```

**Including snapshots of few cause_of_deaths below:**



- **Observation-** The maximum deaths because of **Meningitis** has occurred in India, Nigeria followed by Pakistan. Ethiopia, China, Indonesia, Niger, Democratic Republic of Congo, Mali, Uganda

Countries w/ the Most Alzheimer's Disease and Other Dementias Deaths

- **Observation-** The maximum deaths because of **Alzheimer's Disease and Other Dementias** has occurred in China, United States, Japan, India, Germany, Russia, Italy, Brazil, France, United Kingdom



Countries w/ the Most Malaria Deaths

- **Observation**- The maximum deaths because of **Malaria** has occurred in Nigera, Democratic Republic of Congo, India, Uganda, Burkina faso, Cote d'lvoire, Mozambique, Tanzania, Ghana, Mali

Countries w/ the Most HIV/AIDS Deaths

- **Observation-** The maximum deaths because of **HIV/AIDS** has occurred in Soth Africa, Kenya, Tanzania, India, Nigeria, Uganda, Zimbabwe, Ethiopia, Mozambique, Malawi



Countries w/ the Most Tuberculosis Deaths

- **Observation-** The maximum deaths because of **Tuberculosis** has occurred in India(With a very high count), Indonesia, China, Pakistan, Nigeria, Ethopia, Democratic Republic of Congo, Bangladesh, Myanmar, Philippines

Countries w/ the Most Cardiovascular Diseases Deaths

- **Observation-** The maximum deaths because of **Cardiovascular Diseases** has occurred in China, India, Russia, United States, Indonesia, Ukraine, Germany, Brazil, Japan, Pakistan



Countries w/ the Most Lower Respiratory Infections Deaths

- **Observation-** The maximum deaths because of **Lower Respiratory Infections** has occurred in India, China, Nigeria, Japan, Pakistan, Ethopia, United States, Bangladesh, Brazil, Indonesia

Countries w/ the Most Neonatal Disorders Deaths

- **Observation-** The maximum deaths because of **Neonatal Disorders** has occurred in India(Again withn a very high count), Pakistan, Nigeria , China, Ethiopia, Bangladesh,, Indonesia Brazil, Democratic Republic of Congo, Tanzania



Countries w/ the Most Self-harm Deaths

- **Observation-** The maximum deaths because of **Self-harm** has occurred in India, China, Russia, United States, Japan, Ukraine, France, Pakistan, Germany, Brazil

Countries w/ the Most Diarrheal Diseases Deaths

- **Observation-** The maximum deaths because of **Diarrheal Diseases** has occurred in India, Nigeria, Pakistan, Indonesia, Ethiopia, Democratic Republic of Congo, Bangladesh, Niger, China, Angola


Countries w/ the Most Neoplasms Deaths

- **Observation-** The maximum deaths because of **Neoplasms** has occurred in China, United States, India, Japan, Russia, Germany, Brazil, United Kingdom, Italy, France

Countries w/ the Most Conflict and Terrorism Deaths

- **Observation-** The maximum deaths because of **Conflict and Terrorism** has occurred in Rwanda, Syria, Iraq, Afganisthan, Burundi, Ethopia, Democratic Republic of Congo, Yemen, Sudan, Sri Lanka.



Countries w/ the Most Diabetes Mellitus Deaths

- **Observation-** The maximum deaths because of Diabetes Mellitus has occurred in India, China, United States, Indonesia, Mexico, Brazil, Pakistan, Germany, South Africa, Bangladesh

Countries w/ the Most Road Injuries Deaths

- **Observation-** The maximum deaths because of **Road Injuries** has occurred in China(With a very high count), India, United States, Brazil, Indonesia, Russia, Iran, Egypt, Democratic Republic of Congo, Thailand



Countries w/ the Most Chronic Respiratory Diseases Deaths

- **Observation**- The maximum deaths because of Chronic Respiratory Diseases has occurred in China, India, United States, Indonesia, Pakistan, Bangladesh, Brazil, Russia, United Kingdom, Myanmar

Countries w/ the Most Cirrhosis and Other Chronic Liver Diseases Deaths

- **Observation-** The maximum deaths because of **Cirrhosis and Other Chronic Liver Diseases** has occurred in India, China, Indonesia, United States, Egypt, Russia, Pakistan, Nigeria, Mexico, Brazil



Countries w/ the Most Digestive Diseases Deaths

- **Observation-** The maximum deaths because of **Digestive Diseases** has occurred in India, China, Indonesia, United States, Russia, Brazil, Nigeria, Egypt, Pakistan, Mexico

**Other Key Observations:**

- The other important observation noticed was India is leading in death-rates and that too with so high ratio in most causes-of-death, which shows there is lot of scope of improvement in medical facilities in India.

- The other countries following India in most death rates in most causes-of-death is China, Indonesia, United States, Ethiopia(Even though population is less, death rates comparatively are very high)

# 5. Displaying each Country's total and top 10 causes of death

- For displaying top 10 causes and total count's death, we have selected first 11 largest count entries of each country and displayed in the form of bar graph.

```
In [19]:  # for Country/Territory in the data set
          for x in df_report_mod[0:-1].index:
              # group all the rows by Country/Territory column and grab the 10 highest values
              data = df_report_mod.set_index('Country/Territory').iloc[x].nlargest(11)

              y = df_report_mod['Country/Territory'].iloc[x]

              # configurations for the bar graphs
              plt.figure(figsize=(10, 6))
              plt.bar(data = data, x = data.index, height = data.values, width=.6, color = ['blue','red','green','magenta','yellow','cyan',
              plt.xticks(rotation = 90)
              plt.xlabel('Causes of Death', size= 13)
              plt.ylabel('Total Deaths (1990-2019)')
              plt.title(y + "'s Top 10 Causes of Death")
```
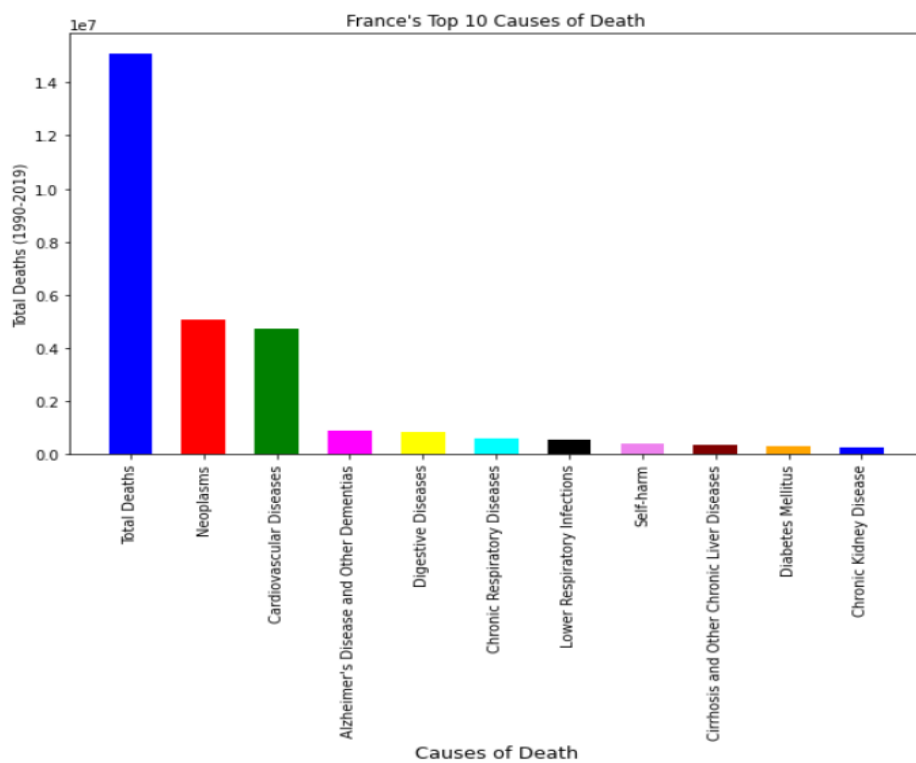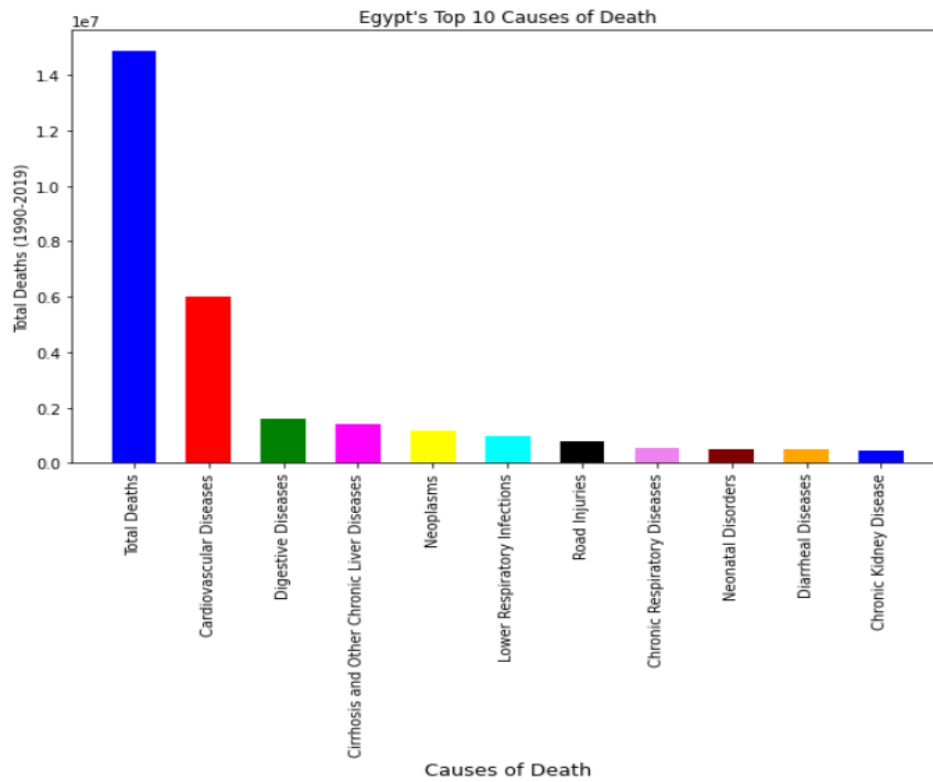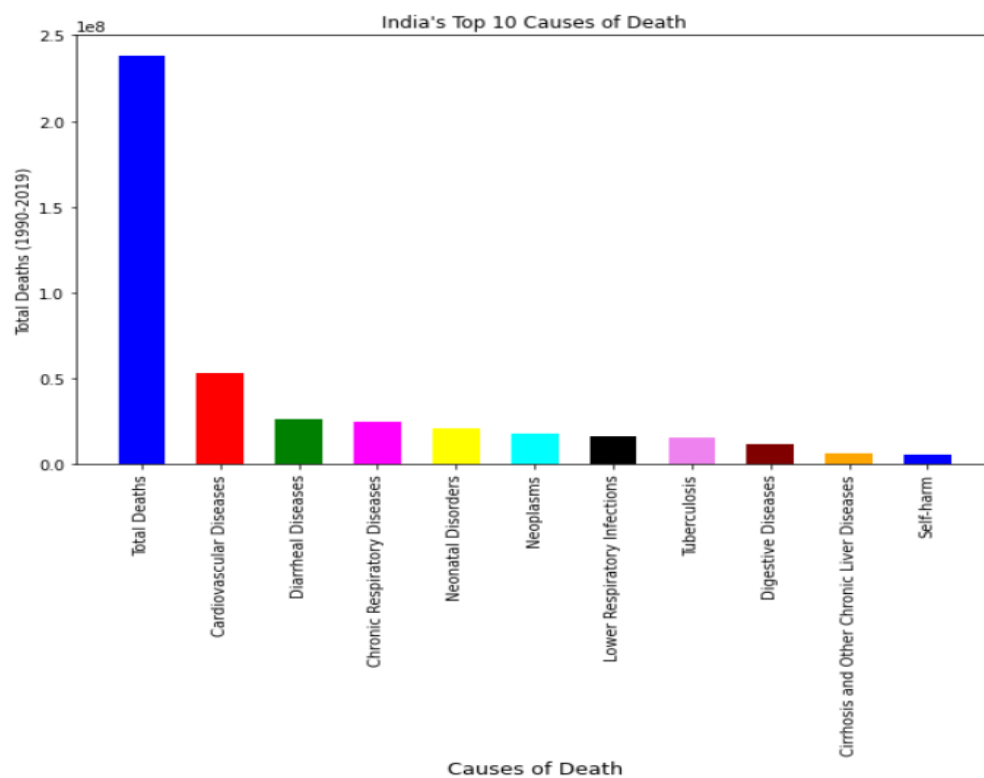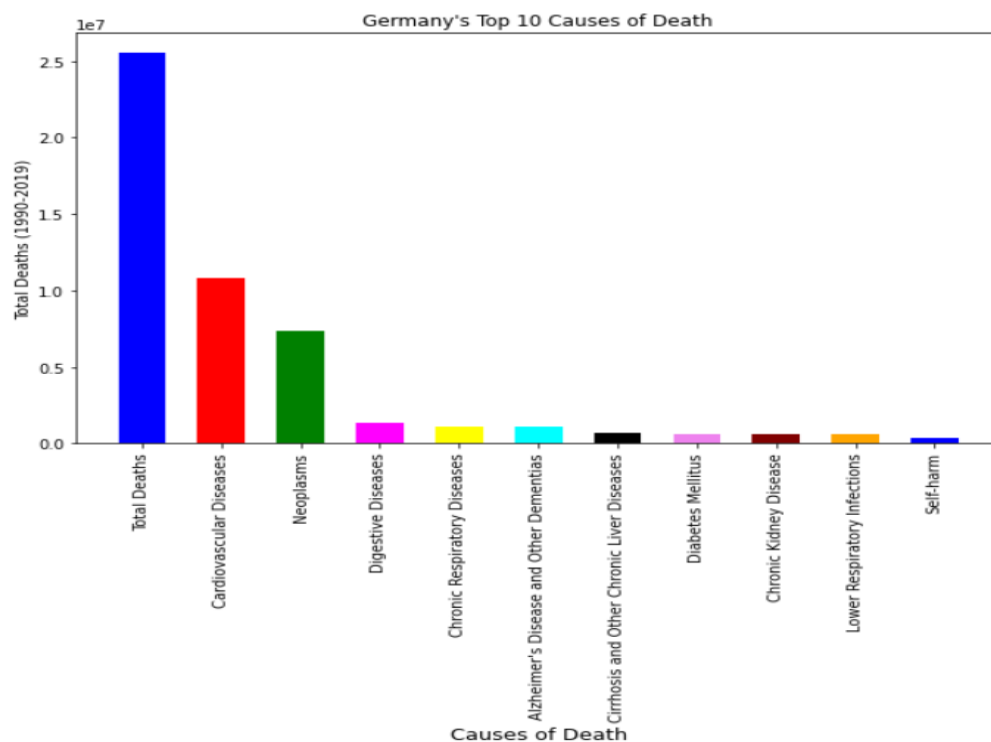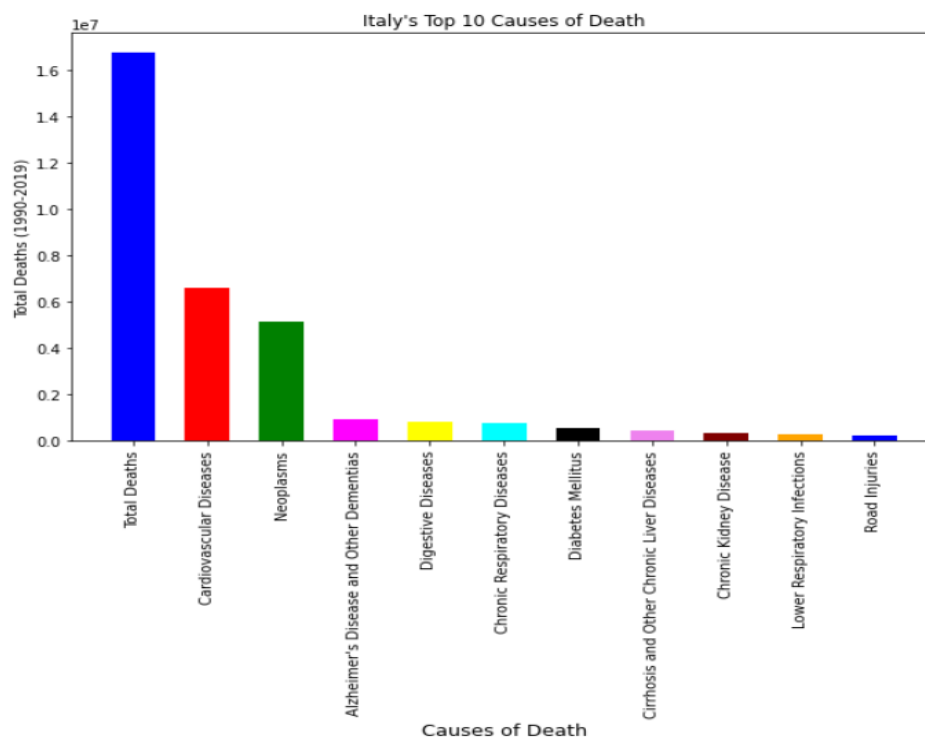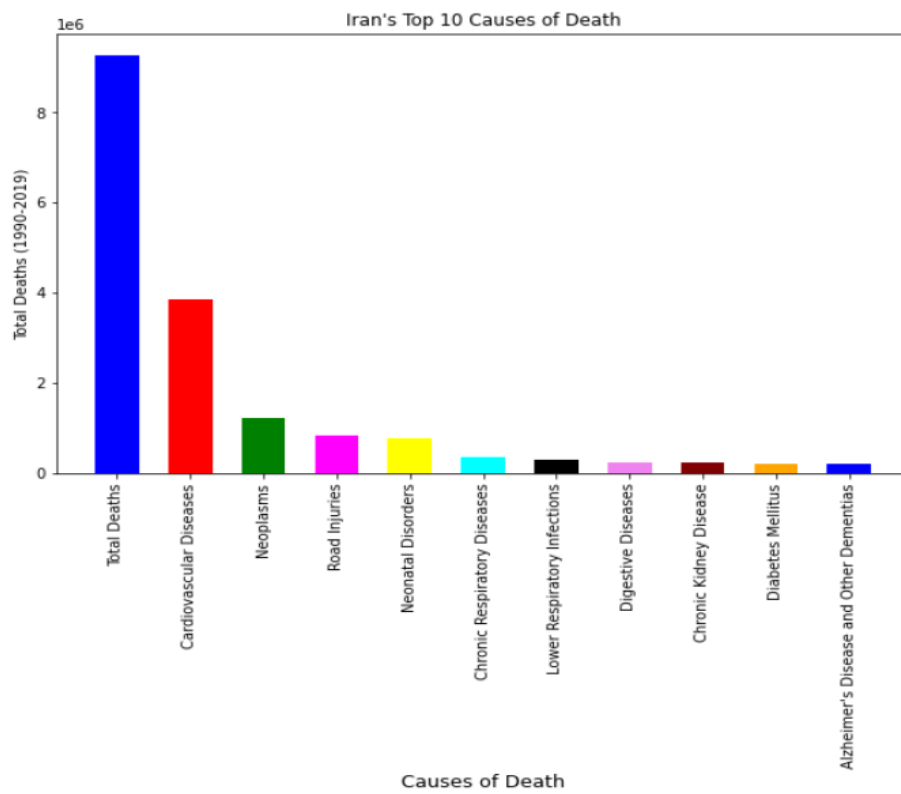
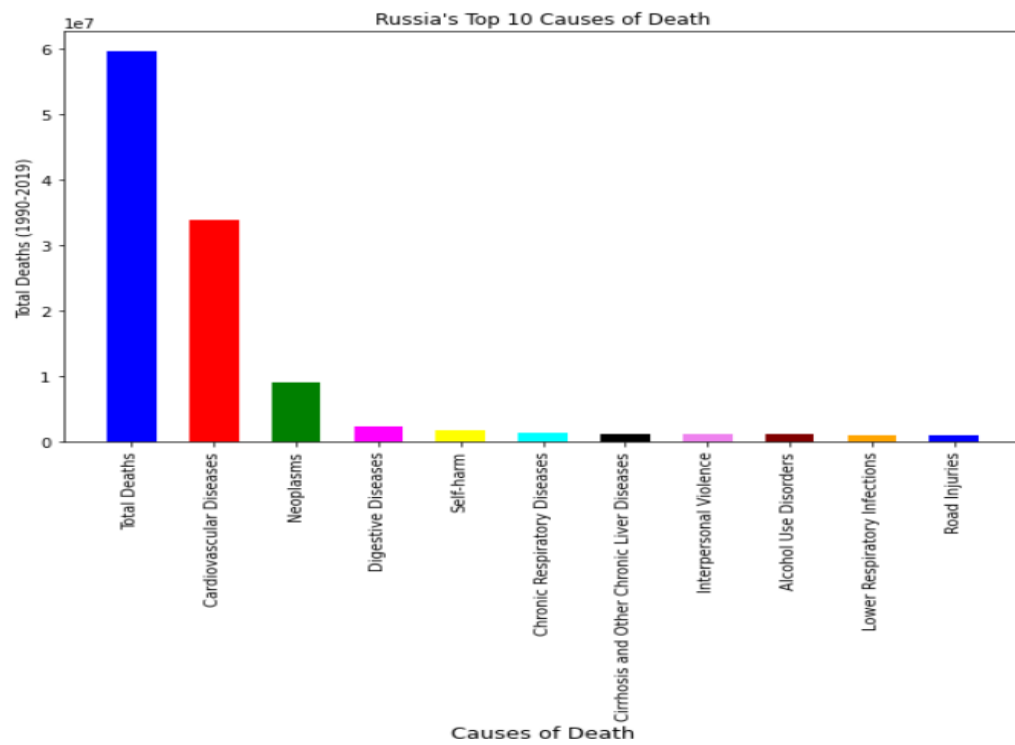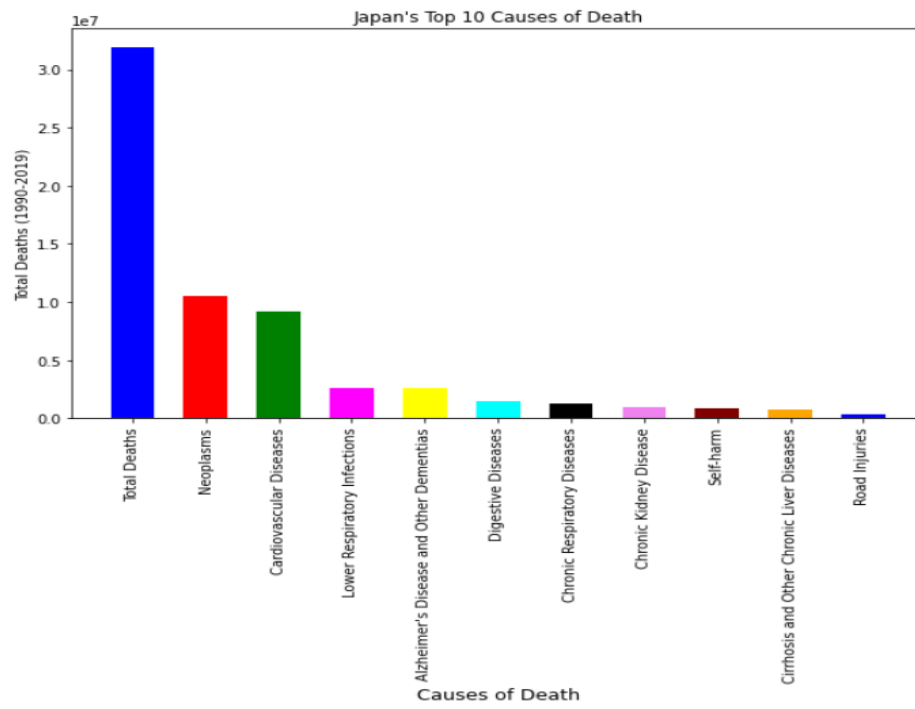**Attaching snapshots of few Countries with their top causes of death below:**

Afghanistan's Top 10 Causes of Death



Australia's Top 10 Causes of Death

Bangladesh's Top 10 Causes of Death



Canada's Top 10 Causes of Death

Egypt's Top 10 Causes of Death

Causes of Death



France's Top 10 Causes of Death

Causes of Death

Germany's Top 10 Causes of Death

India's Top 10 Causes of Death

Iran's Top 10 Causes of Death



Italy's Top 10 Causes of Death

Japan's Top 10 Causes of Death



Russia's Top 10 Causes of Death

South Africa's Top 10 Causes of Death



United Arab Emirates's Top 10 Causes of Death
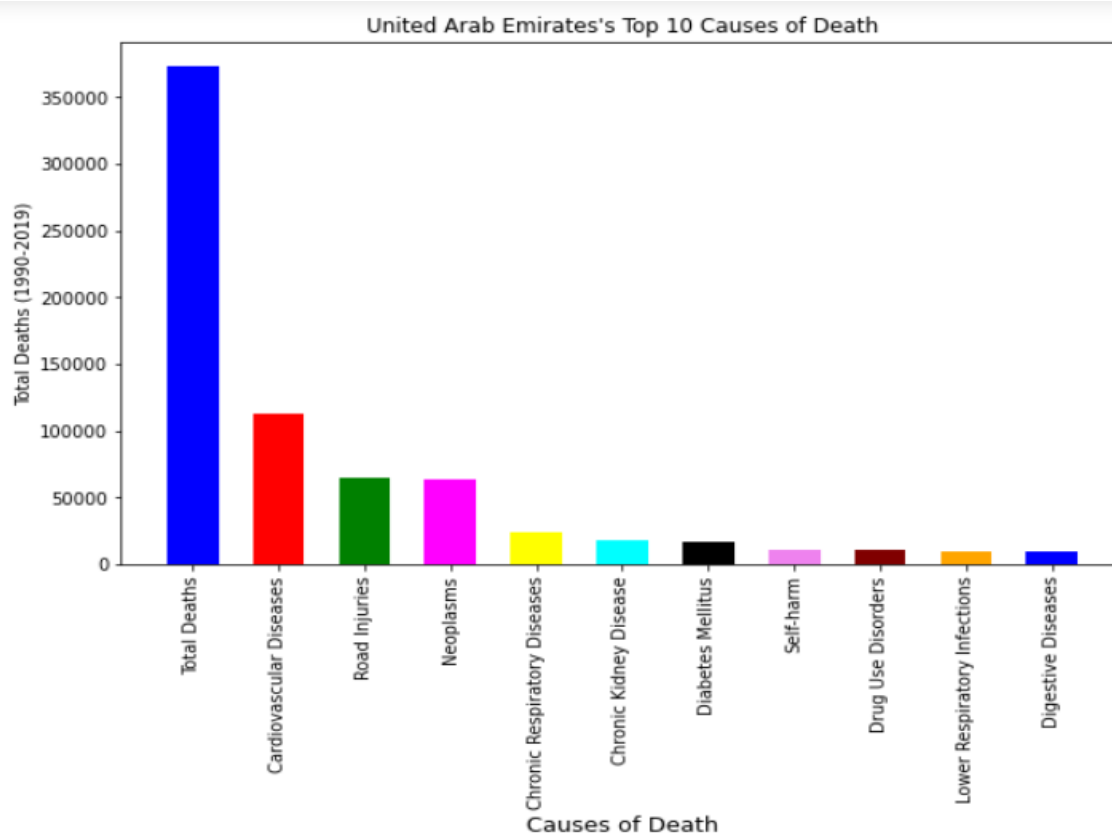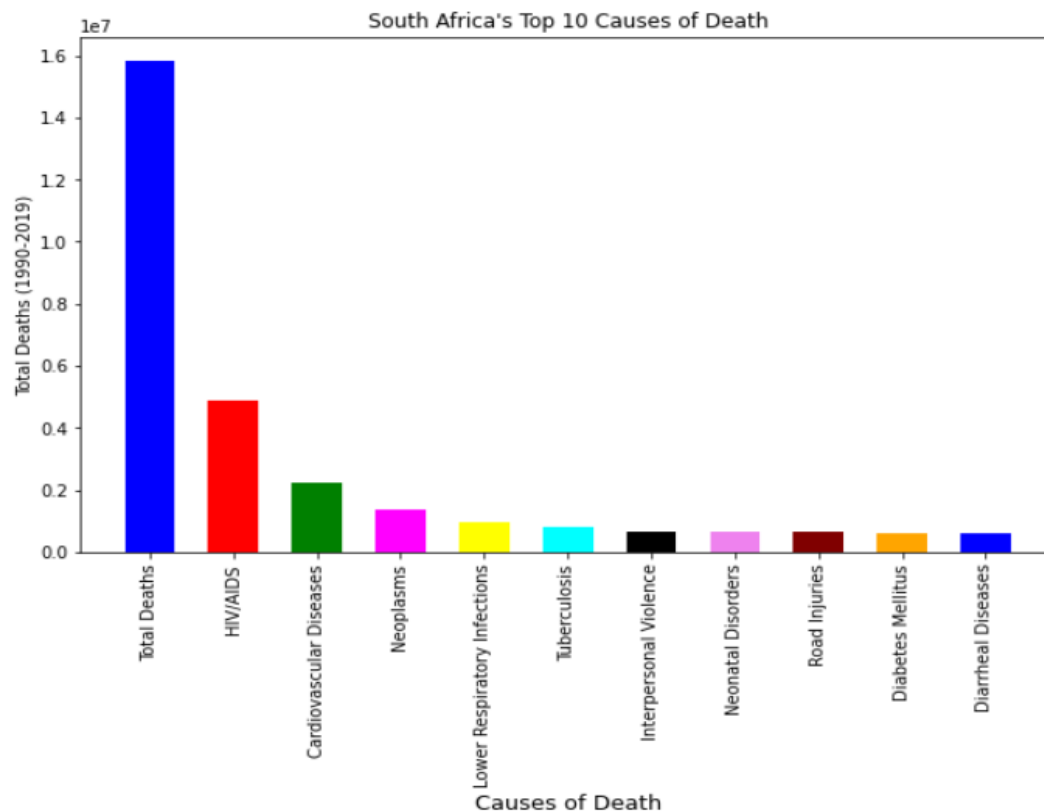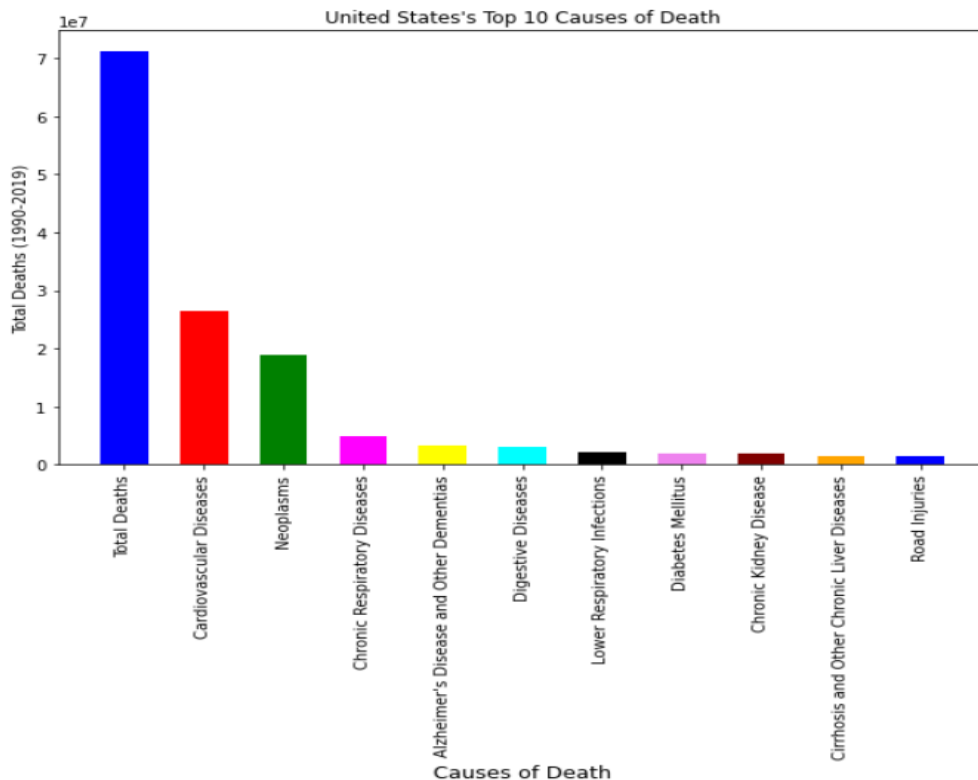
United States's Top 10 Causes of Death

**Key Observations:**

- The main cause of death in most of the countries is "Cardiovascular Diseases".

- The second most cause of death in most countries is "Neoplasms".

- The other leading causes of deaths are "Chronic Respiratory Diseases", "Lower Respiratory Infections", "Digestive Diseases", "Neonatal Disorders", "Diabetes Mellitus

- "Neonatal Disorders" being one of the leading cause of deaths is a deep concern since babies/children dying because of certain problems from birth can be overcome if proper facilities are developed to check for the diseases before birth

# Result & Conclusion

The overall analysis helped in exploring Cause-of- Death of all countries as well as gave Country profile related to the diseases.

Poor healthcare and nutrition are the main reason of deaths in many countries and thus measures for improvement shall be practiced to increase life not just of patients but their family members too. Neonatal disorders too being one of the top cause of death is  a serious concern and shows insufficient facilities for baby or toddlers.

# Limitations of this work and Scope for Future Work

There are many factors that can help us in getting much more detailed insights of mortality rate to improve the models' predictions.

The `biggest limitation of the provided dataset was that there wasn't any population related data, and thus we could not judge the mortality rate in terms of population percentage. Also, specifications related to gender, sex, age-group of population would have helped in extracting more meaningful data.