

Project by : KRITIKA RAWAT .  
project title : Analysis of COVID-19 data .  
course : bechlor's in computer science.

### Project Title: Analysis of Global COVID-19 Data

#### Objective:

The primary goal of this project was to analyze a comprehensive dataset of global COVID-19 cases to uncover trends and patterns in the spread of the virus. By visualizing the data, the project aimed to highlight the impact of the pandemic across different countries and provide actionable insights that could aid in public health planning and resource allocation.

#### Approach:

##### 1. Data Preparation:

- **Data Import:** The dataset, which contained information on COVID-19 cases from various countries, was imported into a Pandas DataFrame.
- **Data Cleaning:** The dataset was cleaned to ensure accuracy and consistency. This included:
  - Handling missing values by either imputing them with appropriate values or removing incomplete records.
  - Converting data types where necessary, such as changing date columns to datetime objects for better manipulation.
  - Dropping irrelevant columns that did not contribute to the analysis, thus focusing only on the necessary data.
- **Preprocessing:** Ensured that all data was in a format suitable for analysis, including standardizing country names and verifying data integrity.

##### 2. Exploratory Data Analysis (EDA):

- **Pattern Identification:** Conducted initial explorations to identify key patterns and trends in the data.
- **Top Countries Analysis:** Focused on visualizing the distribution of confirmed COVID-19 cases. Created an interactive bar chart to highlight the top 10 countries with the

highest number of confirmed cases, making it easier to see which countries were most affected.

### 3. Visualization:

- **Plotly Express Library:** Used Plotly Express to create interactive and dynamic visualizations. This library allows users to explore the data visually and interact with the charts for a more comprehensive understanding.
- **Pie Chart:** Developed an interactive pie chart to show the distribution of critical COVID-19 cases by country, providing a clear view of which countries had the highest number of critical cases.

### Key Insights:

- **High Impact Countries:** The analysis identified the USA, India, and France as having the highest numbers of confirmed COVID-19 cases. This finding underscored the severe impact of the pandemic in these regions, highlighting the need for targeted public health interventions and resource allocation.

### Tools & Technologies:

- **Python:** The primary programming language used for data analysis and visualization.
- **Pandas:** Utilized for data manipulation, including cleaning, preprocessing, and aggregation.
- **Plotly Express:** Employed to create interactive and engaging visualizations, making it easier to understand complex data and trends.

### Impact:

The project provides a detailed, data-driven perspective on the global impact of COVID-19. By visualizing the data effectively, it enables stakeholders, including public health officials and policymakers, to better comprehend the spread of the virus. The insights gained from this analysis can help prioritize resources, plan interventions, and improve strategies for managing public health crises.



ANALIZATION OF COVID-19 DATA.

import pandas as pd

```
# Specify the full path to your CSV file
file_path = "C:/Users/KRITIKA RAWAT/Downloads/covid_dataset.csv"

# Read the CSV file into a DataFrame
df = pd.read_csv(file_path, encoding='unicode_escape')
```

In [7]: df.shape

Out[8]: (225, 9)

In [9]: df.head(25)

Unnamed: 0	country	code	confirmed	recovered	critical	deaths	lastChange	lastUpdate
0	0	Afghanistan	AF	234174	211080	0	7996	2024-06-04T00:16:51+00:00
1	2	Albania	AL	354663	330233	0	3605	2024-06-04T00:27:05+00:00
2	3	Algeria	DZ	272010	183061	0	6881	2024-06-04T00:19:02+00:00
3	5	Andorra	AD	48015	47663	0	165	2024-02-11T22:06:54+00:00
4	6	Angola	AO	107327	103419	0	1937	2024-02-11T22:06:54+00:00
...	...	...	...	...	...	...	...	...
220	243	Wallis and Futuna	WF	3500	438	0	8	2024-02-11T22:03:31+00:00
221	244	Western Sahara	EH	10	9	0	1	2024-02-11T22:11:46+00:00
222	245	Yemen	YE	11945	9124	0	2159	2024-02-11T22:03:34+00:00
223	246	Zambia	ZM	349304	341316	0	4089	2024-06-04T00:26:18+00:00
224	247	Zimbabwe	ZW	286359	258888	12	5740	2024-06-04T00:28:51+00:00

225 rows x 9 columns

In [11]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 225 entries, 0 to 224
Data columns (total 9 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   unnamed:0   225 non-null    int64
 1   country     225 non-null    object
 2   code        224 non-null    object
 3   confirmed   225 non-null    int64
 4   recovered   225 non-null    int64
 5   critical     225 non-null    int64
 6   deaths      225 non-null    int64
 7   lastChange  225 non-null    object
 8   lastUpdate  225 non-null    object
dtypes: int64(6), object(3)
memory usage: 15.9+ MB
```

In [12]: df = df.drop(columns=['Unnamed: 0'])

```
df['lastChange'] = pd.to_datetime(df['lastChange'])
df['lastUpdate'] = pd.to_datetime(df['lastUpdate'])
```

this graph shows Top 10 Countries by Confirmed COVID-19 Cases

In [45]: import plotly.express as px

```
# Create a DataFrame with the top 8 countries
top_countries = df.sort_values(by='confirmed', ascending=False).head(10)

# Create an interactive bar chart
fig = px.bar(top_countries, x='confirmed', y='country',
             title='Top 10 Countries by Confirmed COVID-19 Cases',
             labels={'confirmed': 'Confirmed Cases', 'country': 'Country'},
             text='confirmed')

# Update layout for better readability
fig.update_layout(xaxis_title='Confirmed Cases',
                  yaxis_title='Country',
                  xaxis=dict(type='log',
                           ticks=dict(format='')),
                  title=dict(x=0.5,
                             xanchor='center'))

# Show the plot
fig.show()
```



THE ABOVE SHOWS THE TOP 10 COUNTRIES BY HIGHLY CONFIRMED COVID-19 CASES. USA, INDIA, FRANCE ARE IN TOP 3. MOST HIGH CONFIRMED CASES CAME FROM USA.

In [14]: # Display the first few rows of the DataFrame

print(df.head())

# Get summary statistics

print(df.describe())

# Check for any missing values

print(df.isnull().sum())

```
country code confirmed recovered critical deaths \
0 Afghanistan AF 234174 211080 0 7996
1 Albania AL 354663 330233 0 3605
2 Algeria DZ 272010 183061 0 6881
3 Andorra AD 48015 47663 0 165
4 Angola AO 107327 103419 0 1937

lastChange lastUpdate
0 2024-06-04 00:16:51+00:00 2024-06-09 09:13:41+08:00
1 2024-06-04 00:27:05+00:00 2024-06-09 09:13:41+08:00
2 2024-06-04 00:19:02+00:00 2024-06-09 09:13:41+08:00
3 2024-02-11 22:06:54+00:00 2024-06-09 09:13:41+08:00
4 2024-06-04 00:26:18+00:00 2024-06-09 09:13:41+08:00

confirmed recovered critical deaths
count 2.250000e+02 2.250000e+02 2.250000e+02 2.250000e+02
mean 1.113281e+05 1.043234e+05 1.276667e+01 1.117174e+04
std 1.054384e+05 9.532380e+04 7.238973e+00 1.105680e+05
min 1.188889e+01 0.400000e+00 0.000000e+00 0.000000e+00
50% 1.734848e+04 1.547180e+04 0.000000e+00 1.348889e+02
90% 1.383548e+05 1.406220e+05 0.000000e+00 2.250000e+03
95% 1.384548e+05 1.258432e+05 0.000000e+00 1.038300e+04
max 1.113281e+05 1.098344e+05 940.000000e+00 1.219470e+06

country code
1 1
confirmed
0
recovered
0
critical
0
deaths
0
lastChange
0
lastUpdate
0
dtype: int64
```

In [14]: import pandas as pd

import plotly.express as px

```
# Load the dataset
file_path = "C:/Users/KRITIKA RAWAT/Downloads/covid_dataset.csv"
covid_data = pd.read_csv(file_path)

# Clean the dataset by removing the unnecessary 'Unnamed: 0' column
covid_data = covid_data.drop(columns=['Unnamed: 0'])

# Sort the data by deaths and select the top 10 countries
top_10_deaths = covid_data.sort_values(by='deaths', ascending=False).head(10)

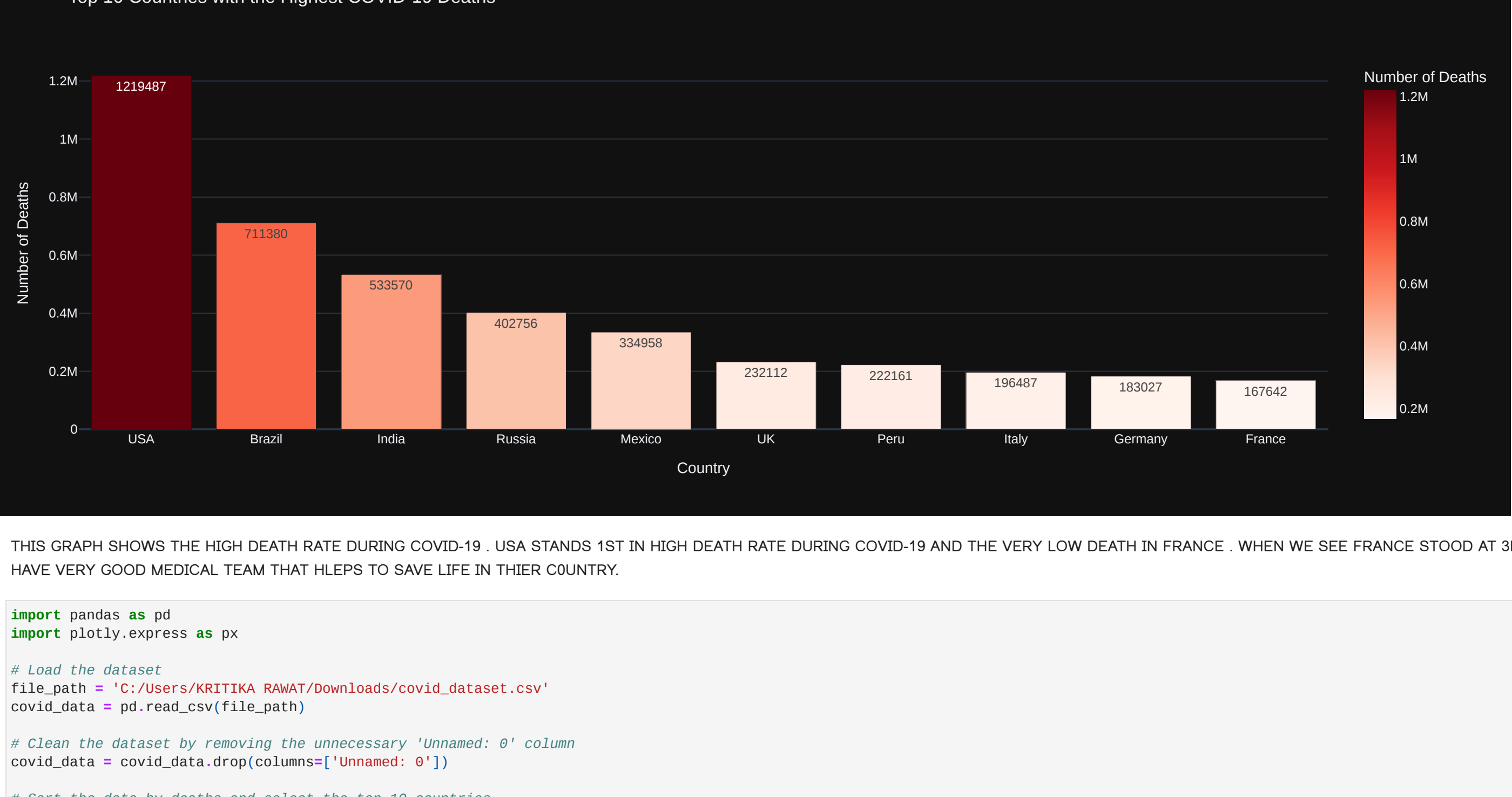
# Create an interactive bar chart using Plotly
fig = px.bar(
    top_10_deaths,
    x='country', # Countries on the x-axis
    y='deaths', # Number of deaths on the y-axis
    text='deaths', # Display the death count on the bars
    title='Top 10 Countries with the Highest COVID-19 Deaths',
    labels={'deaths': 'Number of Deaths', 'country': 'Country'},
    color='deaths', # Color bars based on the number of deaths
    color_continuous_scale='Reds'
)
```

# Update the layout for better visualization

```
fig.update_layout(
    xaxis_title='Country',
    yaxis_title='Number of Deaths',
    template='plotly_dark'
)
```

# Show the interactive bar chart

fig.show()



THIS GRAPH SHOWS THE HIGH DEATH RATE DURING COVID-19. USA STANDS 1ST IN HIGH DEATH RATE DURING COVID-19 AND THE VERY LOW DEATH IN FRANCE. WHEN WE SEE FRANCE STOOD AT 3RD MOST HIGHLY CONFIRMED CASES BUT DUE TO THEY HAVE VERY GOOD MEDICAL TEAM THAT HELPS TO SAVE LIFE IN THEIR COUNTRY.

In [12]: import pandas as pd

import plotly.express as px

```
# Load the dataset
file_path = "C:/Users/KRITIKA RAWAT/Downloads/covid_dataset.csv"
covid_data = pd.read_csv(file_path)
```

# Clean the dataset by removing the unnecessary 'Unnamed: 0' column

covid\_data = covid\_data.drop(columns=['Unnamed: 0'])

# Sort the data by deaths and select the top 10 countries

top\_10\_deaths = covid\_data.sort\_values(by='deaths', ascending=False).head(10)

# Merge the dataset to have a variable column for recovered and deaths

merged\_data = top\_10\_deaths.melt(id\_vars=['country'], value\_vars=['recovered', 'deaths'],

var\_name='case\_type', value\_name='count')

# Create an interactive grouped bar chart using Plotly

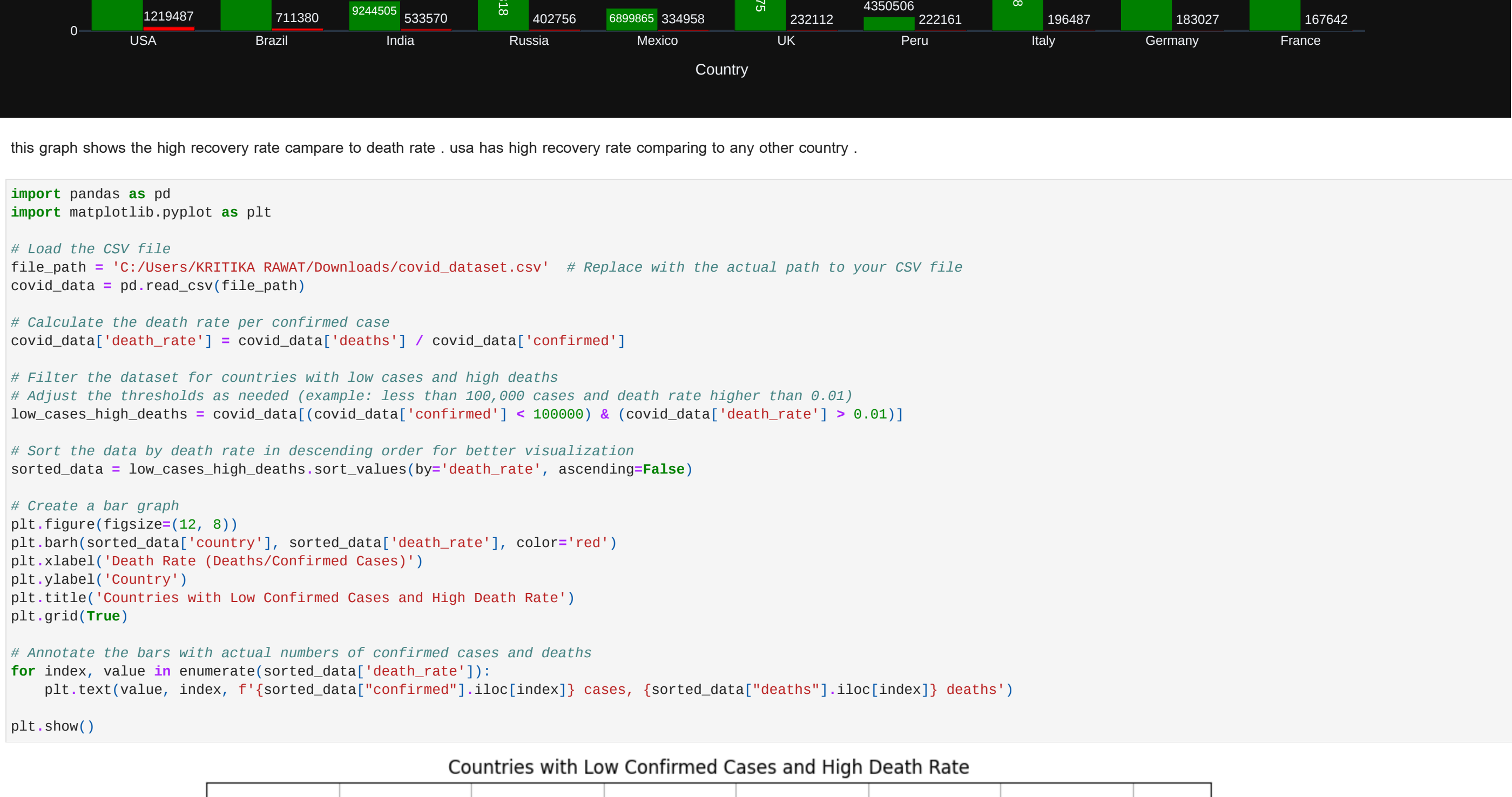
```
fig = px.bar(
    merged_data,
    x='country', # Countries on the x-axis
    y='count', # Number of cases on the y-axis
    color='case_type', # Differentiate between recovered and deaths
    barmode='group', # Group bars side-by-side
    text='count', # Display the count on the bars
    title='Top 10 Countries with the Highest COVID-19 Recoveries and Deaths',
    labels={'count': 'Number of Cases', 'country': 'Country', 'case_type': 'Case Type'},
    color_discrete_map={'recovered': 'green', 'deaths': 'red'} # Set specific colors for recovered and deaths
)
```

# Update the layout for better visualization

```
fig.update_layout(
    xaxis_title='Country',
    yaxis_title='Number of Cases',
    template='plotly_dark'
)
```

# Show the interactive bar chart

fig.show()



this graph shows the high recovery rate compare to death rate . usa has high recovery rate comparing to any other country .

In [18]: import pandas as pd

import matplotlib.pyplot as plt

```
# Load the CSV file
file_path = "C:/Users/KRITIKA RAWAT/Downloads/covid_dataset.csv" # Replace with the actual path to your CSV file
covid_data = pd.read_csv(file_path)
```

# Calculate the death rate per confirmed case

covid\_data['death\_rate'] = covid\_data['deaths'] / covid\_data['confirmed']

# Filter the dataset for countries with low cases and high deaths

# Adjust the thresholds as needed (example: less than 200,000 cases and death rate higher than 0.8)

low\_cases\_high\_deaths = covid\_data[(covid\_data['confirmed'] < 200000) & (covid\_data['death\_rate'] > 0.8)]

# Sort the data by death rate in descending order for better visualization

sorted\_data = low\_cases\_high\_deaths.sort\_values(by='death\_rate', ascending=False)

# Create a bar graph

plt.figure(figsize=(12, 8))

plt.bar(sorted\_data['country'], sorted\_data['death\_rate'], color='red')

plt.xlabel('Death Rate (Deaths/Confirmed Cases)')

plt.ylabel('Country')

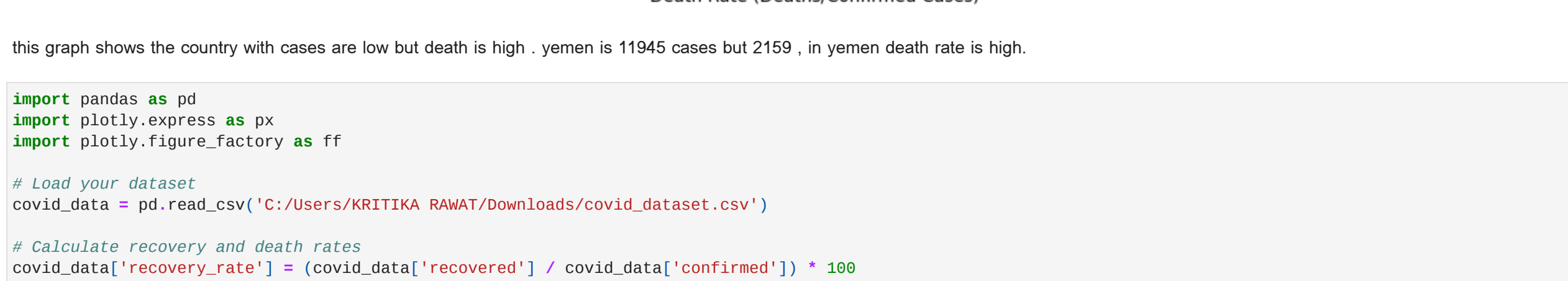
plt.title('Countries with Low Confirmed Cases and High Death Rate')

plt.show()

# Annotate the bars with actual numbers of confirmed cases and deaths

```
for index, value in enumerate(sorted_data['death_rate']):
    plt.text(value, index, f'{sorted_data["confirmed"].iloc[index]} cases, {sorted_data["deaths"].iloc[index]} deaths')

plt.show()
```



this graph shows the country with cases are low but death is high. yemen is 11945 cases but 2159, in yemen death rate is high.

In [17]: import pandas as pd

import plotly.express as px

import plotly.figure\_factory as ff

# Load your dataset

covid\_data = pd.read\_csv('C:/Users/KRITIKA RAWAT/Downloads/covid\_dataset.csv')

# Calculate recovery and death rates

covid\_data['recovery\_rate'] = (covid\_data['recovered'] / covid\_data['confirmed']) \* 100

covid\_data['death\_rate'] = (covid\_data['deaths'] / covid\_data['confirmed']) \* 100

# Scatter plot for Confirmed Cases with Country Name

fig1 = px.scatter(covid\_data, x='confirmed', y='country', size='confirmed', color='confirmed',

hover\_name='country', title='Distribution of Confirmed COVID-19 Cases by Country',

labels={'confirmed': 'Confirmed Cases'})

fig1.show()

# 2. Scatter plot for Recovery Rate with Country Names

fig2 = px.scatter(covid\_data, x='recovery\_rate', y='country', size='recovery\_rate', color='recovery\_rate',

hover\_name='country', title='Recovery Rate by Country',

labels={'recovery\_rate': 'Recovery Rate (%)'})

fig2.show()

# 3. Scatter plot for Death Rate with Country Names

fig3 = px.scatter(covid\_data, x='death\_rate', y='country', size='death\_rate', color='death\_rate',

hover\_name='country', title='Death Rate by Country',

labels={'death\_rate': 'Death Rate (%)'})

fig3.show()

# 4. Scatter plot for Critical Cases with Country Names

fig4 = px.scatter(covid\_data, x='critical', y='country', size='critical', color='critical',

hover\_name='country', title='Distribution of Critical COVID-19 Cases by Country',

labels={'critical': 'Critical Cases'})

fig4.show()

# 5. Correlation Matrix

corr\_matrix = covid\_data[['confirmed', 'recovered', 'critical', 'deaths']].corr()

fig5 = ff.create\_annotated\_heatmap(corr\_matrix.values,

zlist(corr\_matrix.columns),

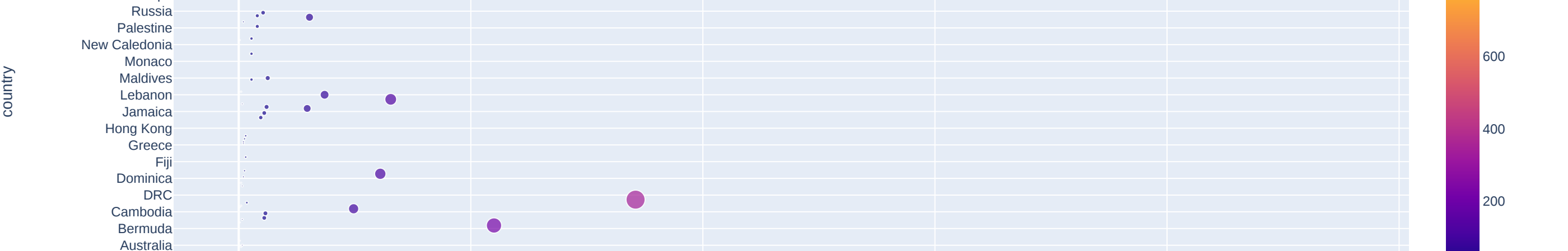
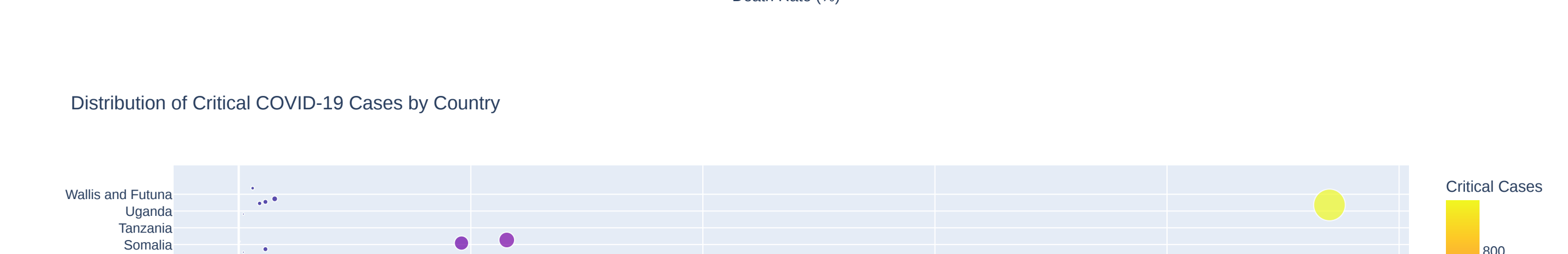
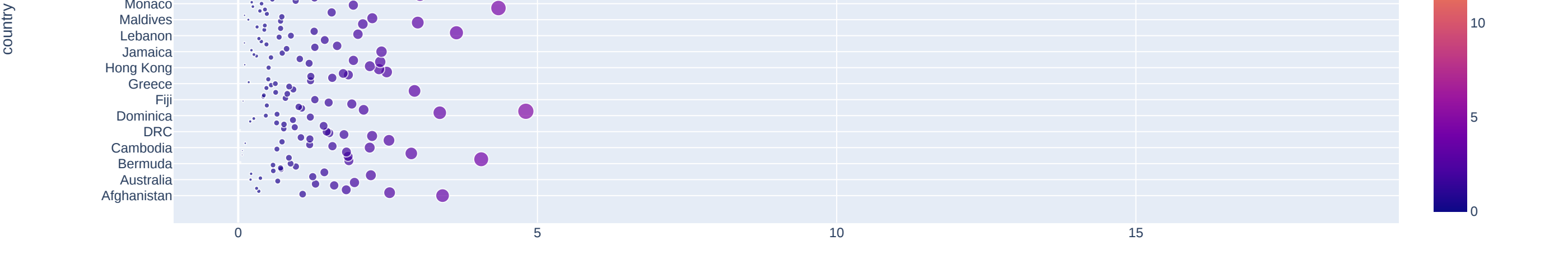
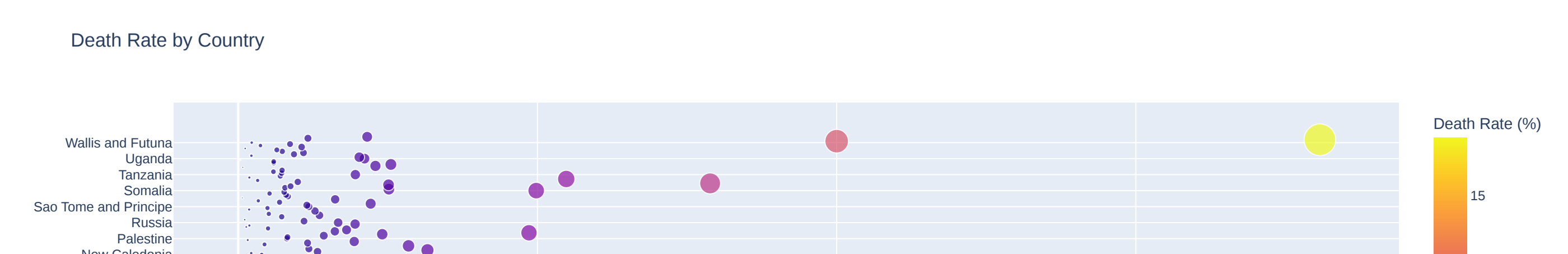
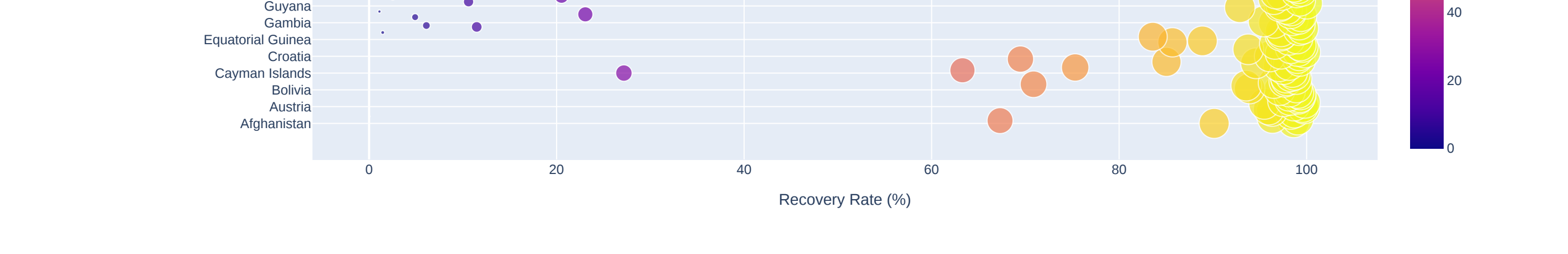
ylist(corr\_matrix.index),

annotation\_text=corr\_matrix.round(2).values,

colorscale='viridis', showscale=True)

fig5.update\_layout(title='Correlation Matrix of COVID-19 Statistics')

fig5.show()



this pie chart shows top 10 countries recorded critical cases



