

DSBDA MINI PROJECT

Ananalysis and visualization of global CO2 emissions

```
1 from google.colab import drive
2 drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 %matplotlib inline
```

DATA PREPROCESSING

```
1 try:
2     df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/DSBDA/CO2 emission by countries.csv", encoding='utf-8')
3 except UnicodeDecodeError:
4     try:
5         df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/DSBDA/CO2 emission by countries.csv", encoding='latin1')
6     except UnicodeDecodeError:
7         df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/DSBDA/CO2 emission by countries.csv", encoding='ISO-8859-1')
```

```
1 df.head()
```

	Country	Code	Calling Code	Year	CO2 emission (Tons)	Population(2022)	Area	% of World	Density(km2)
0	Afghanistan	AF	93	1750	0.0	41128771.0	652230.0	0.40%	63/km²
1	Afghanistan	AF	93	1751	0.0	41128771.0	652230.0	0.40%	63/km²
2	Afghanistan	AF	93	1752	0.0	41128771.0	652230.0	0.40%	63/km²
3	Afghanistan	AF	93	1753	0.0	41128771.0	652230.0	0.40%	63/km²
4	Afghanistan	AF	93	1754	0.0	41128771.0	652230.0	0.40%	63/km²

```
1 missing_values = df.isnull().sum()
2 print("Missing values:\n", missing_values)
```

```
Missing values:
Country          0
Code             2168
Calling Code     3523
Year             0
CO2 emission (Tons)  0
Population(2022)  6504
Area             4336
% of World       4336
Density(km2)     6504
dtype: int64
```

```
1 print("Data types:\n", df.dtypes)
```

```
Data types:
Country          object
Code             object
Calling Code     object
Year            int64
CO2 emission (Tons) float64
Population(2022) float64
Area            float64
% of World       object
Density(km2)     object
dtype: object
```

```
1 print("Summary statistics:\n", df.describe())
```

Summary statistics:		Year	CO2 emission (Tons)	Population(2022)	Area
count	59620.000000		5.962000e+04	5.311600e+04	5.528400e+04
mean	1885.000000		1.034774e+09	3.992260e+07	6.522073e+05
std	78.231085		1.041652e+10	1.482365e+08	1.865483e+06
min	1750.000000		0.000000e+00	1.131200e+04	2.100000e+01
25%	1817.000000		0.000000e+00	1.770414e+06	1.770450e+04
50%	1885.000000		0.000000e+00	8.673095e+06	1.103815e+05
75%	1953.000000		8.715092e+06	2.862920e+07	4.925730e+05
max	2020.000000		4.170000e+11	1.425887e+09	1.709824e+07

```
1 df.isnull().sum()
```

```
Country          0
Code             2168
Calling Code     3523
Year             0
CO2 emission (Tons)  0
Population(2022)  6504
Area             4336
```

```
% of World      4336
Density(km2)    6504
dtype: int64
```

```
1 df.shape

(59620, 9)
```

```
1 # Drop rows with missing values
2 df.dropna(inplace=True)
```

```
1 df.isnull().sum()

Country      0
Code         0
Calling Code  0
Year         0
CO2 emission (Tons)  0
Population(2022)  0
Area         0
% of World   0
Density(km2)  0
dtype: int64
```

```
1 df.describe()
```

	Year	CO2 emission (Tons)	Population(2022)	Area
count	48509.000000	4.850900e+04	4.850900e+04	4.850900e+04
mean	1885.000000	1.160748e+09	4.186199e+07	6.328071e+05
std	78.231235	1.125903e+10	1.545682e+08	1.551563e+06
min	1750.000000	0.000000e+00	1.131200e+04	2.100000e+01
25%	1817.000000	0.000000e+00	2.305825e+06	2.633800e+04
50%	1885.000000	0.000000e+00	9.038309e+06	1.303730e+05
75%	1953.000000	1.660158e+07	3.054758e+07	5.516950e+05
max	2020.000000	4.170000e+11	1.425887e+09	9.984670e+06

```
1 df.shape

(48509, 9)
```

```
1 #Fill the Null values
2 df['Population(2022)'].fillna(df['Population(2022)'].mean(), inplace=True)
3 df['Area'].fillna(df['Area'].mean(), inplace=True)
```

✓ LABEL ENCODING & DATA TRANSFORMATION

```
1 # Select columns of a specific data type
2 numeric_columns = df.select_dtypes(include=['float64', 'int64'])
3 categorical_columns = df.select_dtypes(include=['object'])
```

```
1 numeric_columns
```

	Year	CO2 emission (Tons)	Population(2022)	Area
0	1750	0.0	41128771.0	652230.0
1	1751	0.0	41128771.0	652230.0
2	1752	0.0	41128771.0	652230.0
3	1753	0.0	41128771.0	652230.0
4	1754	0.0	41128771.0	652230.0
...
59615	2016	736467042.0	16320537.0	390757.0
59616	2017	746048675.0	16320537.0	390757.0
59617	2018	757903042.0	16320537.0	390757.0
59618	2019	768852126.0	16320537.0	390757.0
59619	2020	779383468.0	16320537.0	390757.0

48509 rows × 4 columns

```
1 categorical_columns
```

	Country	Code	Calling Code	% of World	Density(km2)
0	Afghanistan	AF	93	0.40%	63/km²
1	Afghanistan	AF	93	0.40%	63/km²
2	Afghanistan	AF	93	0.40%	63/km²
3	Afghanistan	AF	93	0.40%	63/km²
4	Afghanistan	AF	93	0.40%	63/km²
...
59615	Zimbabwe	ZW	263	0.30%	42/km²
59616	Zimbabwe	ZW	263	0.30%	42/km²
59617	Zimbabwe	ZW	263	0.30%	42/km²
59618	Zimbabwe	ZW	263	0.30%	42/km²
59619	Zimbabwe	ZW	263	0.30%	42/km²

48509 rows × 5 columns

```
1 from sklearn.preprocessing import LabelEncoder

1 label_encoder = LabelEncoder()

1 df['Country_Encoded'] = label_encoder.fit_transform(df['Country'])

1 df.head()
```

	Country	Code	Calling Code	Year	C02 emission (Tons)	Population(2022)	Area	% of World	Density(km2)	Country_Encoded
0	Afghanistan	AF	93	1750	0.0	41128771.0	652230.0	0.40%	63/km²	0
1	Afghanistan	AF	93	1751	0.0	41128771.0	652230.0	0.40%	63/km²	0
2	Afghanistan	AF	93	1752	0.0	41128771.0	652230.0	0.40%	63/km²	0
3	Afghanistan	AF	93	1753	0.0	41128771.0	652230.0	0.40%	63/km²	0
4	Afghanistan	AF	93	1754	0.0	41128771.0	652230.0	0.40%	63/km²	0

```
1 unique_sum = df['Country'].nunique()
2 print("Sum of unique values in 'Country' column:", unique_sum)

Sum of unique values in 'Country' column: 179
```

```
1 # Drop the 'Code' column from the DataFrame
2 df.drop(columns=['Code'], inplace=True)

1 df.head()
```

	Country	Calling Code	Year	C02 emission (Tons)	Population(2022)	Area	% of World	Density(km2)	Country_Encoded
0	Afghanistan	93	1750	0.0	41128771.0	652230.0	0.40%	63/km²	0
1	Afghanistan	93	1751	0.0	41128771.0	652230.0	0.40%	63/km²	0
2	Afghanistan	93	1752	0.0	41128771.0	652230.0	0.40%	63/km²	0
3	Afghanistan	93	1753	0.0	41128771.0	652230.0	0.40%	63/km²	0
4	Afghanistan	93	1754	0.0	41128771.0	652230.0	0.40%	63/km²	0

```
1 df['Calling Code'] = pd.to_numeric(df['Calling Code'], errors='coerce')
2 df['% of World'] = pd.to_numeric(df['% of World'].str.rstrip('%'), errors='coerce')
3 df['Density(km2)'] = pd.to_numeric(df['Density(km2)'].str.replace('/km²', ''), errors='coerce')

1 df.head()
```

	Country	Calling Code	Year	C02 emission (Tons)	Population(2022)	Area	% of World	Density(km2)	Country_Encoded
0	Afghanistan	93.0	1750	0.0	41128771.0	652230.0	0.4	63.0	0
1	Afghanistan	93.0	1751	0.0	41128771.0	652230.0	0.4	63.0	0
2	Afghanistan	93.0	1752	0.0	41128771.0	652230.0	0.4	63.0	0
3	Afghanistan	93.0	1753	0.0	41128771.0	652230.0	0.4	63.0	0
4	Afghanistan	93.0	1754	0.0	41128771.0	652230.0	0.4	63.0	0

```
1 df.dtypes

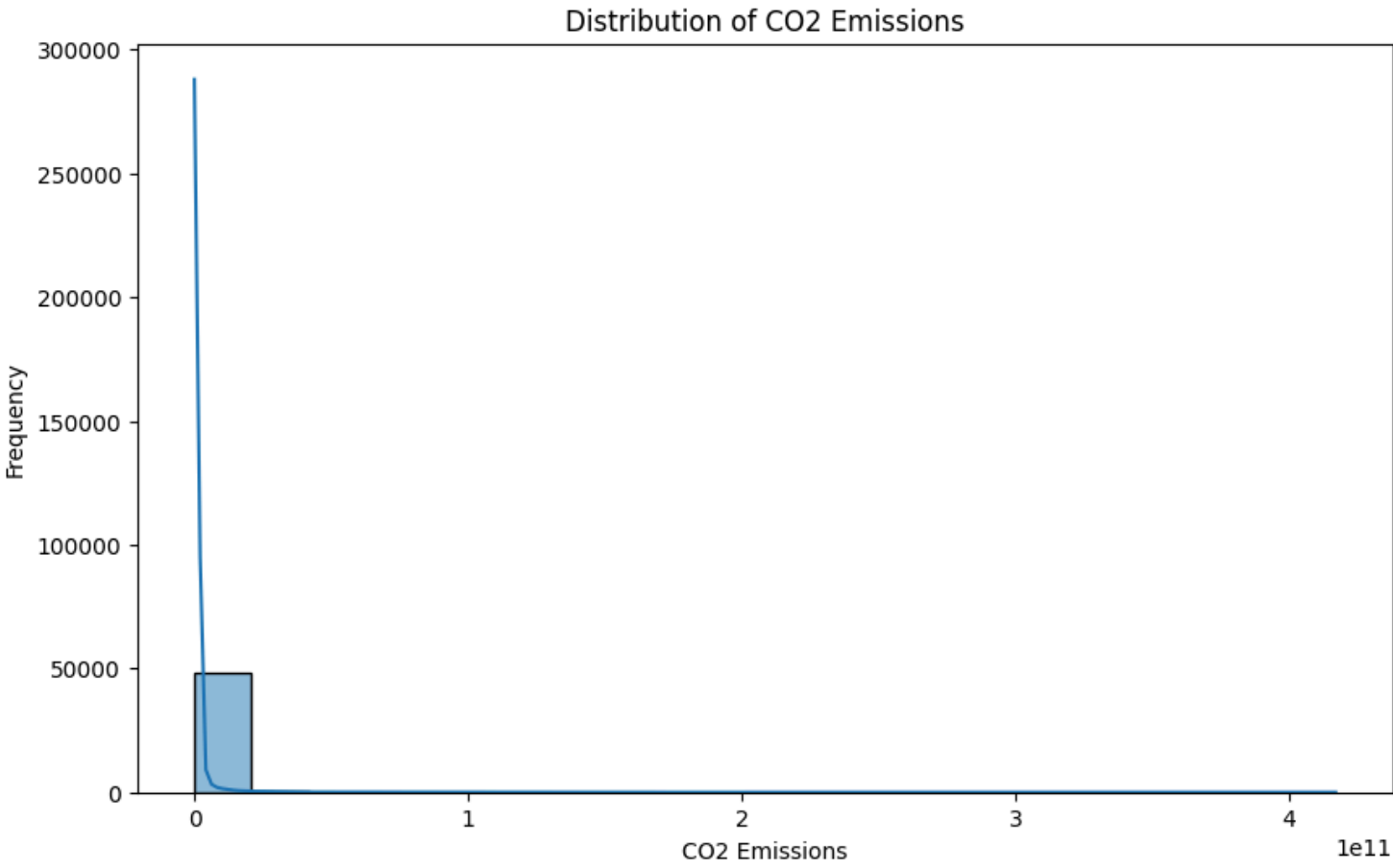
Country          object
Calling Code     float64
Year             int64
C02 emission (Tons) float64
Population(2022) float64
Area             float64
% of World       float64
Density(km2)     float64
```

Country_Encoded int64
dtype: object

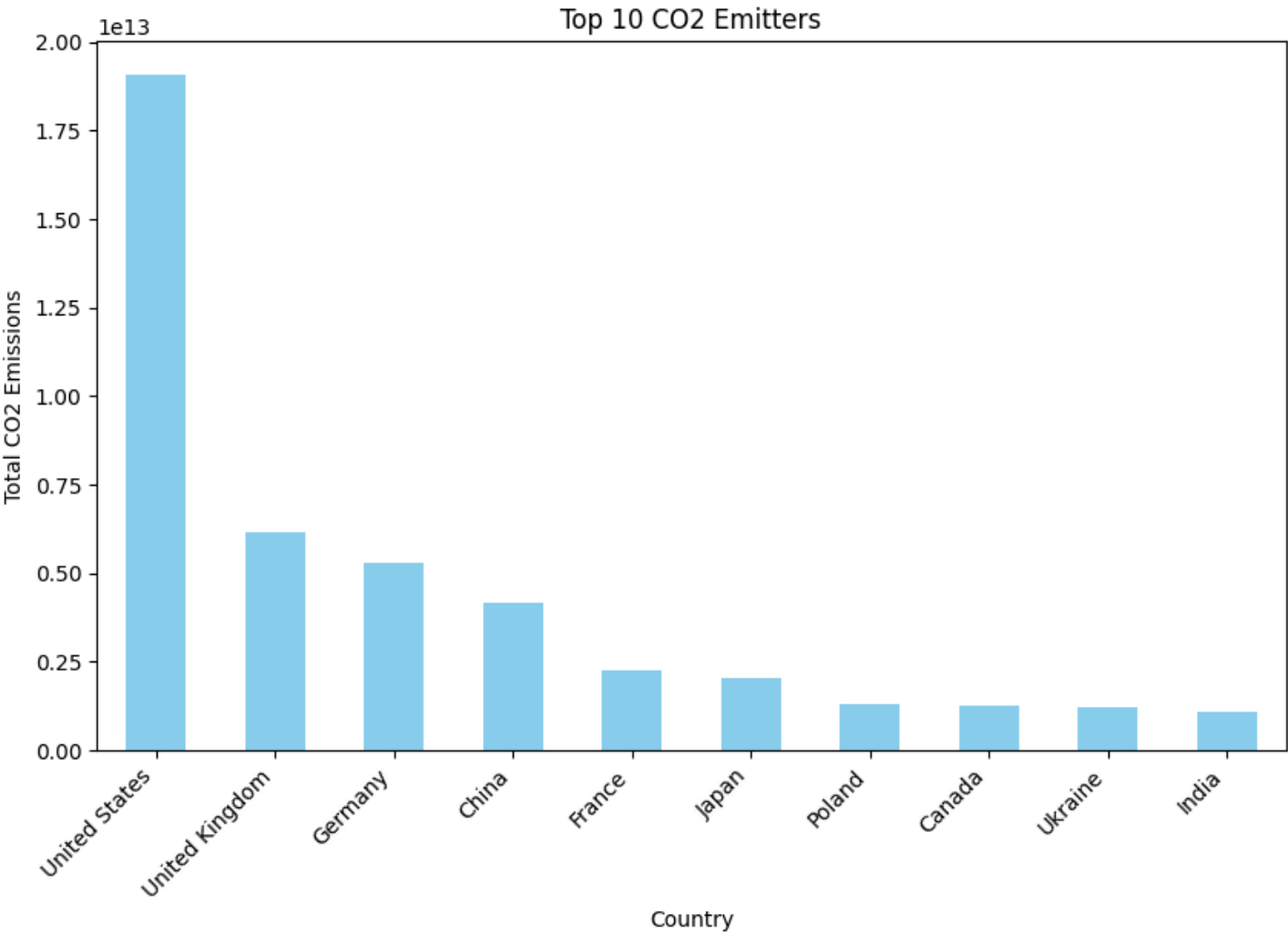
1

DATA VISUALIZATION

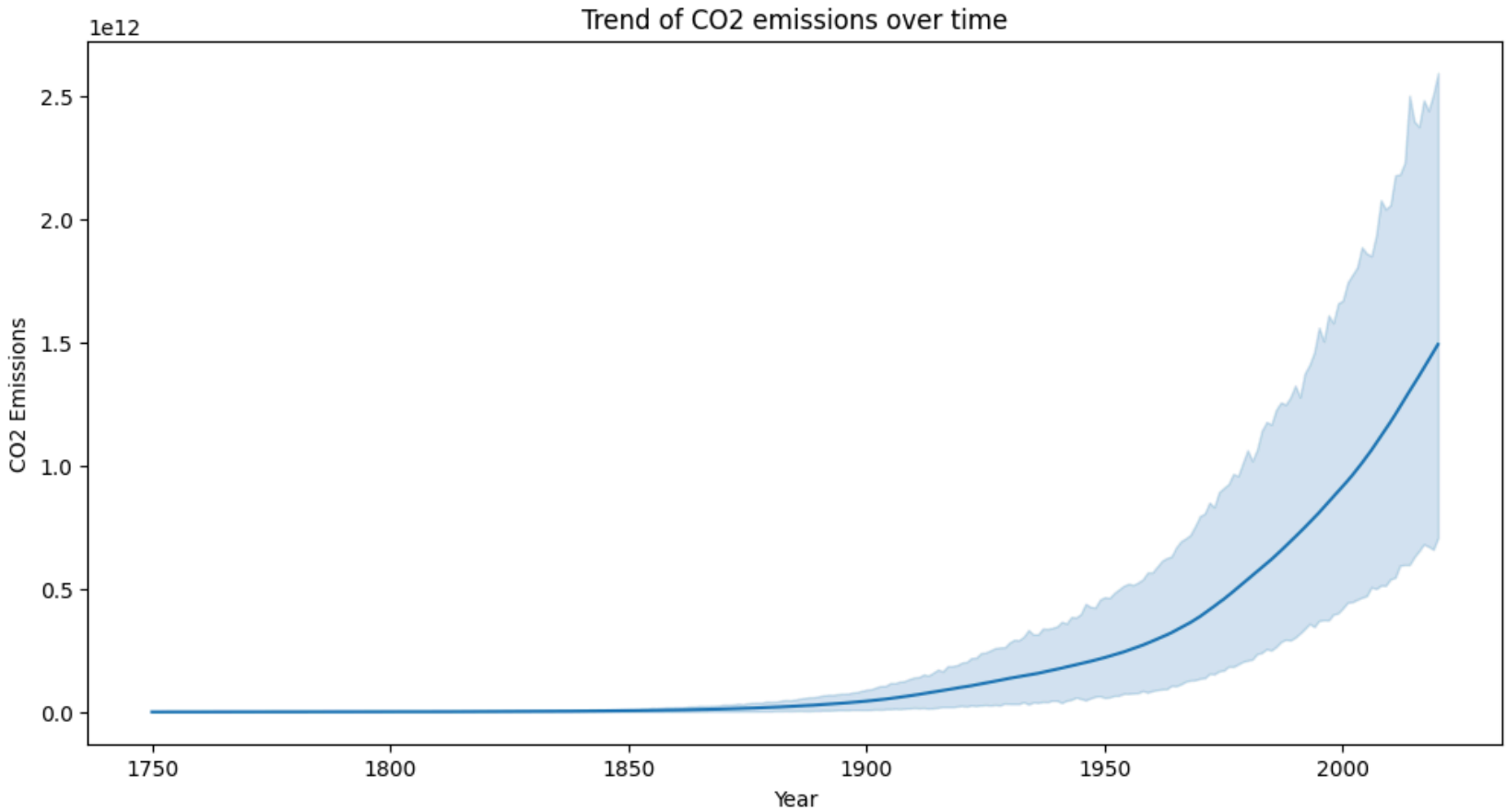
```
1 #Bar plot to visualize top 10 emitters:  
2 plt.figure(figsize=(10, 6))  
3 sns.histplot(df['CO2 emission (Tons)'], bins=20, kde=True)  
4 plt.title('Distribution of CO2 Emissions')  
5 plt.xlabel('CO2 Emissions')  
6 plt.ylabel('Frequency')  
7 plt.show()  
8
```



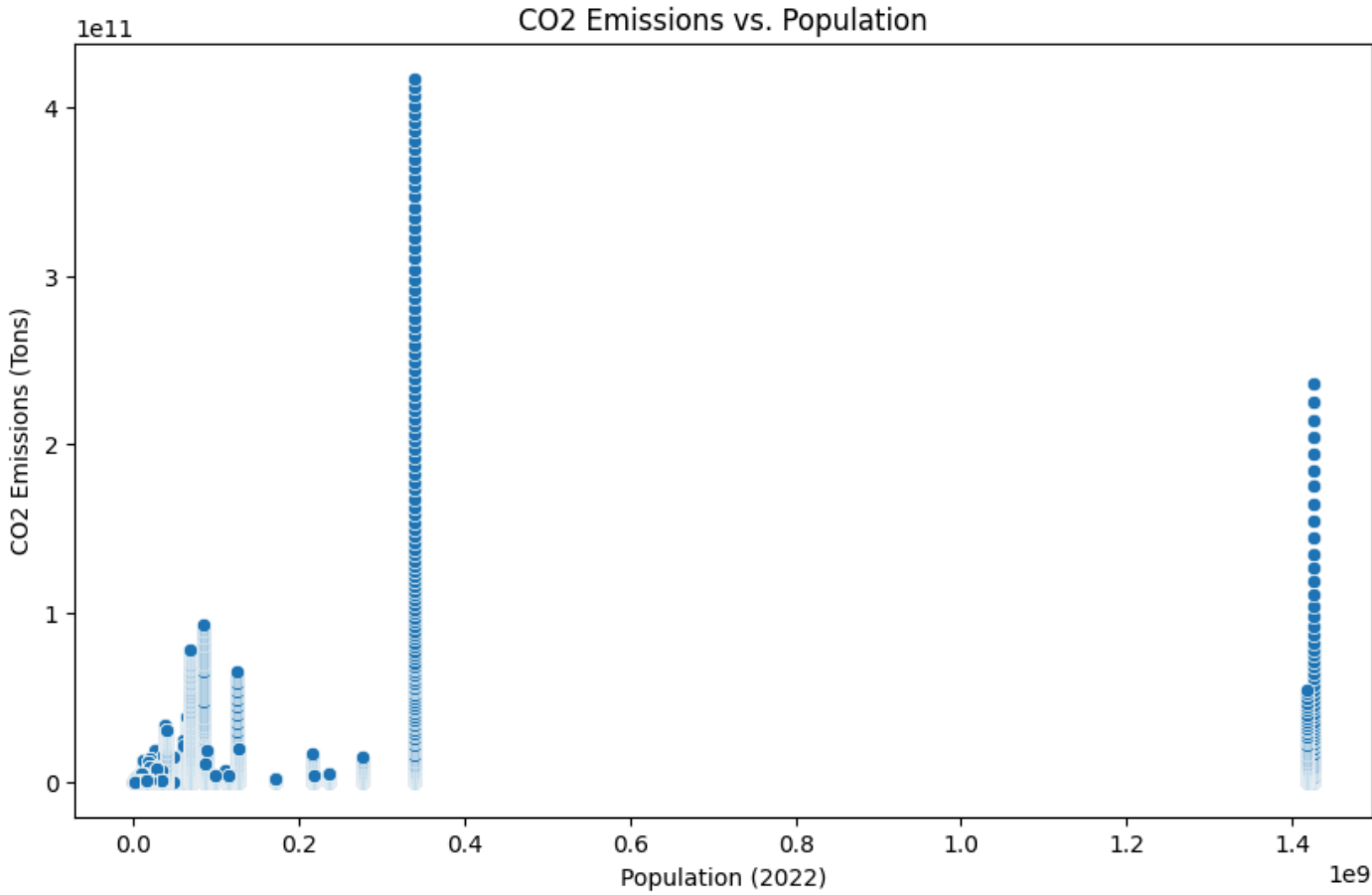
```
1 #Line plot to visualize trend of CO2 emissions over years:  
2 plt.figure(figsize=(10, 6))  
3 top_10_emitters = df.groupby('Country')['CO2 emission (Tons)'].sum().nlargest(10)  
4 top_10_emitters.plot(kind='bar', color='skyblue')  
5 plt.title('Top 10 CO2 Emitters')  
6 plt.xlabel('Country')  
7 plt.ylabel('Total CO2 Emissions')  
8 plt.xticks(rotation=45, ha='right')  
9 plt.show()  
10
```



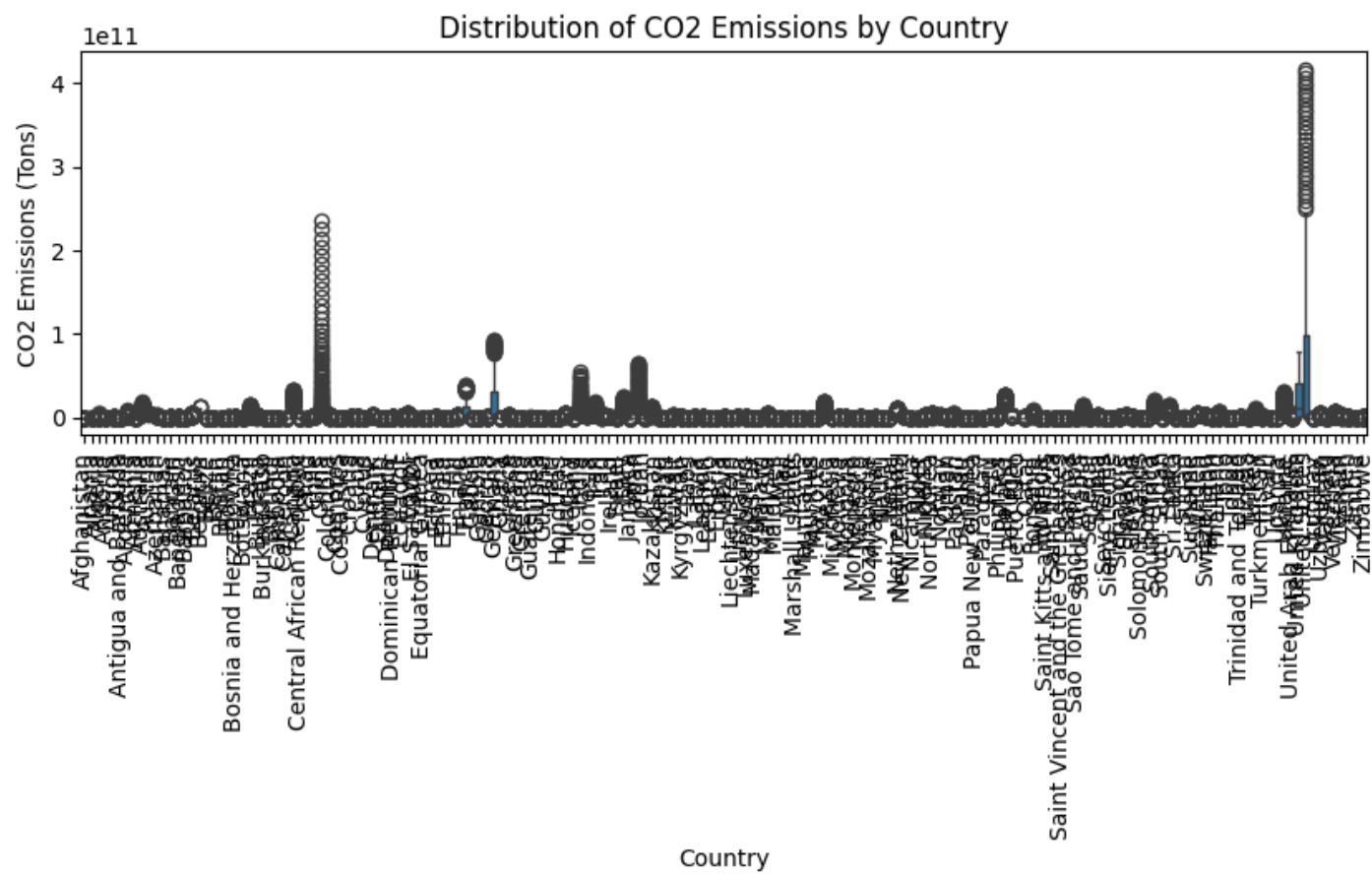
```
1 #Scatter plot to visualize relationship between C02 emissions and population:
2 plt.figure(figsize=(12, 6))
3 sns.lineplot(data=df, x='Year', y='C02 emission (Tons)', estimator='sum')
4 plt.title('Trend of C02 emissions over time')
5 plt.xlabel('Year')
6 plt.ylabel('C02 Emissions')
7 plt.show()
8
```



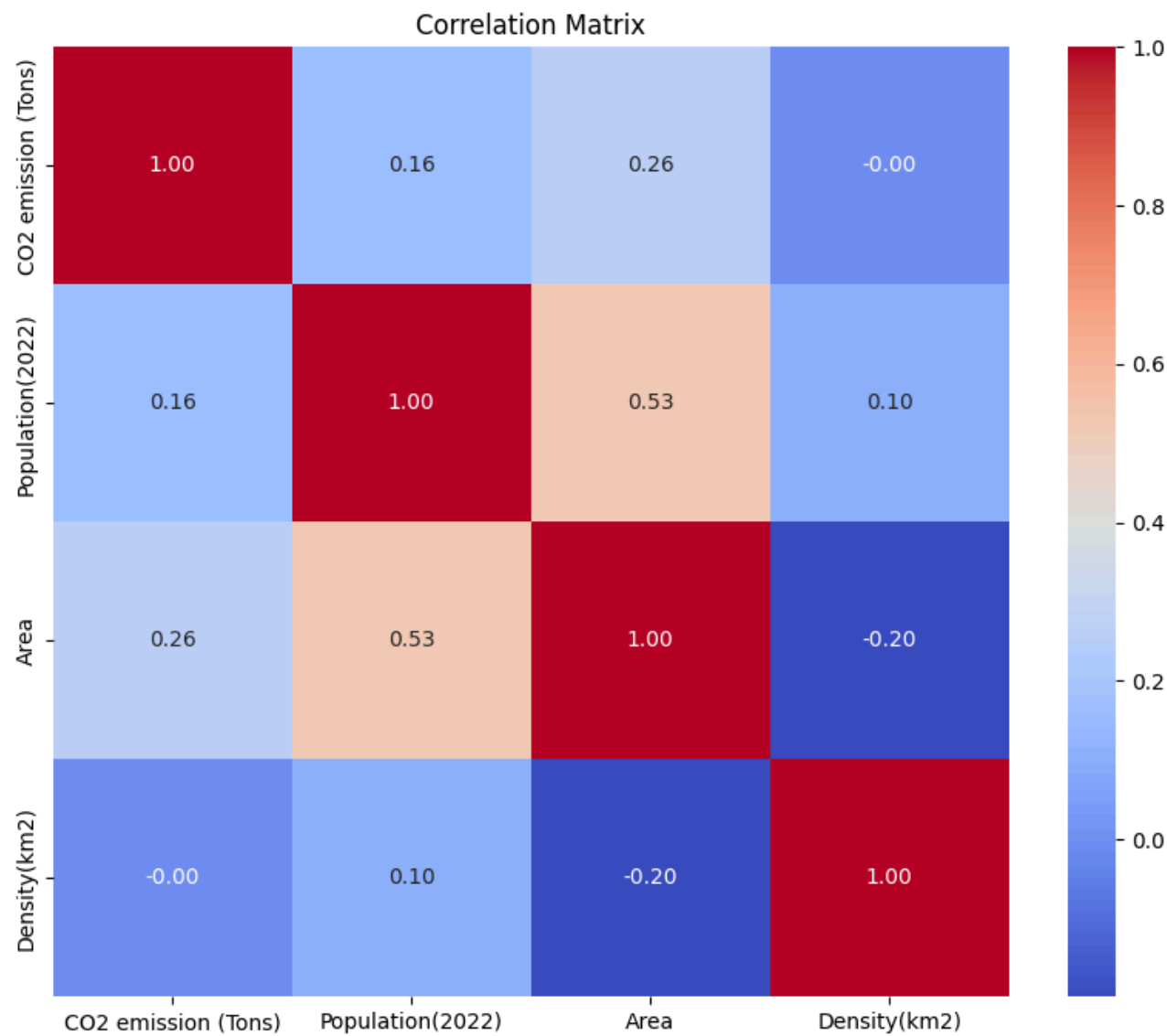
```
1 #Box plot to visualize distribution of C02 emissions by country:
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(data=df, x='Population(2022)', y='C02 emission (Tons)')
4 plt.title('C02 Emissions vs. Population')
5 plt.xlabel('Population (2022)')
6 plt.ylabel('C02 Emissions (Tons)')
7 plt.show()
8
```



```
1 #Box plot to visualize distribution of C02 emissions by country:
2 plt.figure(figsize=(10, 3))
3 sns.boxplot(data=df, x='Country', y='C02 emission (Tons)')
4 plt.title('Distribution of C02 Emissions by Country')
5 plt.xlabel('Country')
6 plt.ylabel('C02 Emissions (Tons)')
7 plt.xticks(rotation=90)
8 plt.show()
9
```



```
1 # Calculate correlation matrix
2 correlation_matrix = df[['CO2 emission (Tons)', 'Population(2022)', 'Area', 'Density(km2)']].corr()
3
4 # Visualize correlation matrix as a heatmap
5 plt.figure(figsize=(10, 8))
6 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
7 plt.title('Correlation Matrix')
8 plt.show()
9
```



```
1 # Total CO2 emissions by country(First 5)
2 total_emissions_by_country = df.groupby('Country')['CO2 emission (Tons)'].sum()
3 print(total_emissions_by_country[:5])
4
5 # Average CO2 emissions by year(First 5)
6 average_emissions_by_year = df.groupby('Year')['CO2 emission (Tons)'].mean()
7 print(average_emissions_by_year[:5])
8
9 # Total CO2 emissions by year(First 5)
10 total_emissions_by_year = df.groupby('Year')['CO2 emission (Tons)'].sum()
11 print(total_emissions_by_year[:5])
12
```

Country	
Afghanistan	3.754106e+09
Albania	9.250246e+09
Algeria	9.582024e+10
Andorra	2.360722e+08
Angola	1.143158e+10

```
Name: C02 emission (Tons), dtype: float64
Year
1750      52237.586592
1751     104475.173184
1752     156733.229050
1753     208991.284916
1754     261269.810056
Name: C02 emission (Tons), dtype: float64
Year
1750      9350528.0
1751     18701056.0
1752     28055248.0
1753     37409440.0
1754     46767296.0
Name: C02 emission (Tons), dtype: float64
```

```
1 # Total C02 emissions by country
2 total_emissions_by_country = df.groupby("Country")["C02 emission (Tons)"].sum().sort_values(ascending=False)
3
4 # Average C02 emissions per capita by country
5 df["C02_emission_per_capita"] = df["C02 emission (Tons)"] / df["Population(2022)"]
6 average_emissions_per_capita = df.groupby("Country")["C02_emission_per_capita"].mean().sort_values(ascending=False)
7 print(average_emissions_per_capita)
```

```
Country
United Kingdom    336.819255
Belgium           250.196083
Germany           234.607068
United States     208.000489
Estonia           176.211194
...
Niger              0.130255
Ethiopia           0.127935
Chad               0.101320
Burundi            0.077897
Puerto Rico       0.023932
Name: C02_emission_per_capita, Length: 179, dtype: float64
```

▼ BUILDING PREDICTION MODEL

```
1 !pip install scikit-learn

Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.2.2)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.25.2)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.11.4)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.4.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.4.0)
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error
4 import matplotlib.pyplot as plt
```

```
1 print(df.head())

   Country  Calling Code  Year  C02 emission (Tons)  Population(2022) \
0  Afghanistan         93.0  1750              0.0         41128771.0
1  Afghanistan         93.0  1751              0.0         41128771.0
2  Afghanistan         93.0  1752              0.0         41128771.0
3  Afghanistan         93.0  1753              0.0         41128771.0
4  Afghanistan         93.0  1754              0.0         41128771.0

   Area  % of World  Density(km2)  Country_Encoded  \
0  652230.0        0.4          63.0              0
1  652230.0        0.4          63.0              0
2  652230.0        0.4          63.0              0
3  652230.0        0.4          63.0              0
4  652230.0        0.4          63.0              0

   C02_emission_per_capita
0              0.0
1              0.0
2              0.0
3              0.0
4              0.0
```

```
1 print(df['Country_Encoded'].nunique())

179
```

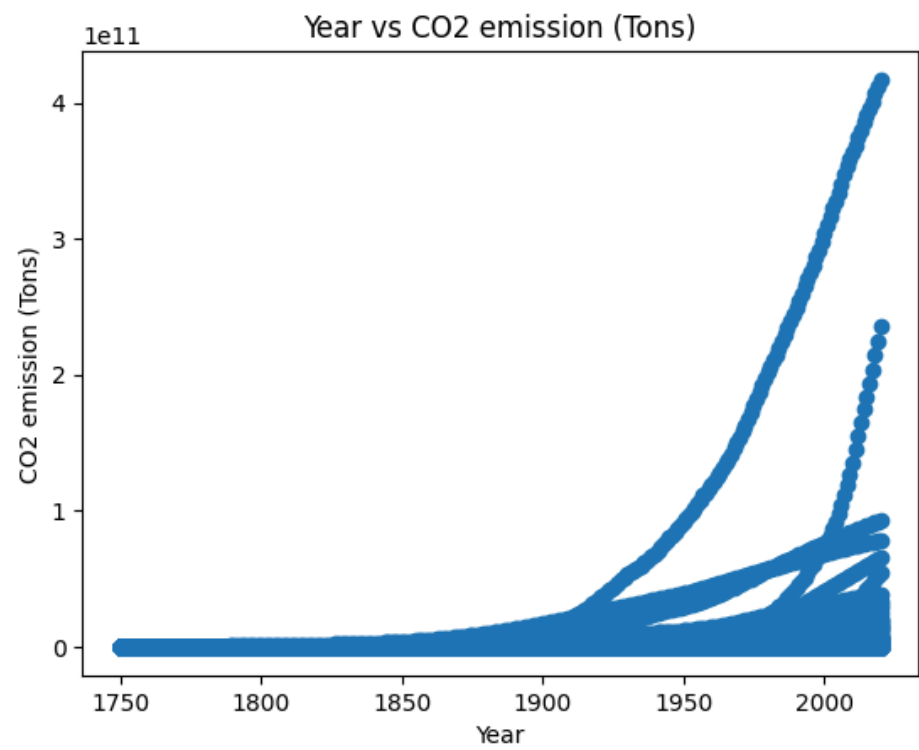
```
1 print(numeric_columns.corr())

           Year  C02 emission (Tons)  Population(2022) \
Year          1.000000e+00          0.134150 -7.233739e-15
C02 emission (Tons)  1.341501e-01          1.000000  1.611472e-01
Population(2022)    -7.233739e-15          0.161147  1.000000e+00
Area               -1.120522e-14          0.262216  5.254822e-01

           Area
Year          -1.120522e-14
C02 emission (Tons)  2.622159e-01
Population(2022)    5.254822e-01
Area              1.000000e+00
```

```
1 X = df[["Year", "Population(2022)"]]
2 y = df["CO2 emission (Tons)"]

1 # Plot Year against CO2 emission (Tons)
2 plt.scatter(X['Year'], y)
3 plt.xlabel('Year')
4 plt.ylabel('CO2 emission (Tons)')
5 plt.title('Year vs CO2 emission (Tons)')
6 plt.show()
```



```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
2
3 # Train a linear regression model
4 model = LinearRegression()
5 model.fit(X_train, y_train)
```

▼ LinearRegression

LinearRegression()

```
1 y_pred = model.predict(X_test)
```

```
1 y_pred

array([-6.18290346e+08,  7.61704583e+08,  2.57100939e+09, ...,
        3.60970282e+09,  2.77131983e+08,  6.82118759e+08])
```

```
1 mse = mean_squared_error(y_test, y_pred)
2 print("Mean Squared Error:", mse)
```

Mean Squared Error: 1.923505065747443e+20