

Data Analysis and Visualization using Python

Alip Yalikun

Virginia Tech

CS5764 Information Visualization

web link <https://dashapp-473886336048.us-east1.run.app/>

Table of Content

Table of Figures	Page 3 - 22
Abstract	Page 23
Introduction	Page 24
Description of Dataset	Page 24
Data pre-processing	Page 24
Outlier detection and removal	Page 24
PCA	Page 25
Normality test	Page 25
Data transformation	Page 25
Heatmap	Page 25
Statistics	Page 25
Data Visualization	Page 26
Subplot	Page 26
Dashboard	Page 26-32

Observation	Page 32-34
Conclusion	Page 34

Table of Figures

Figure 1.

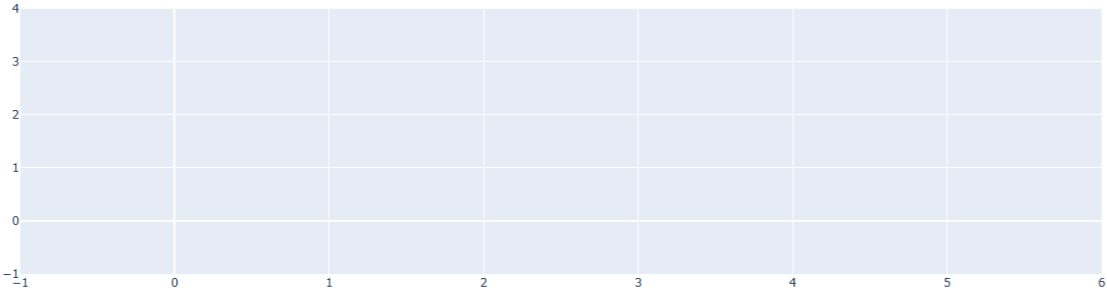
Data cleaning

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

Data Cleaning Methods

Dataset contains 15778 duplicate rows and 0 missing values.

CLEAN DATA



Data Cleaning Methods

All duplicates and missing values have been removed. Dataset is now clean.

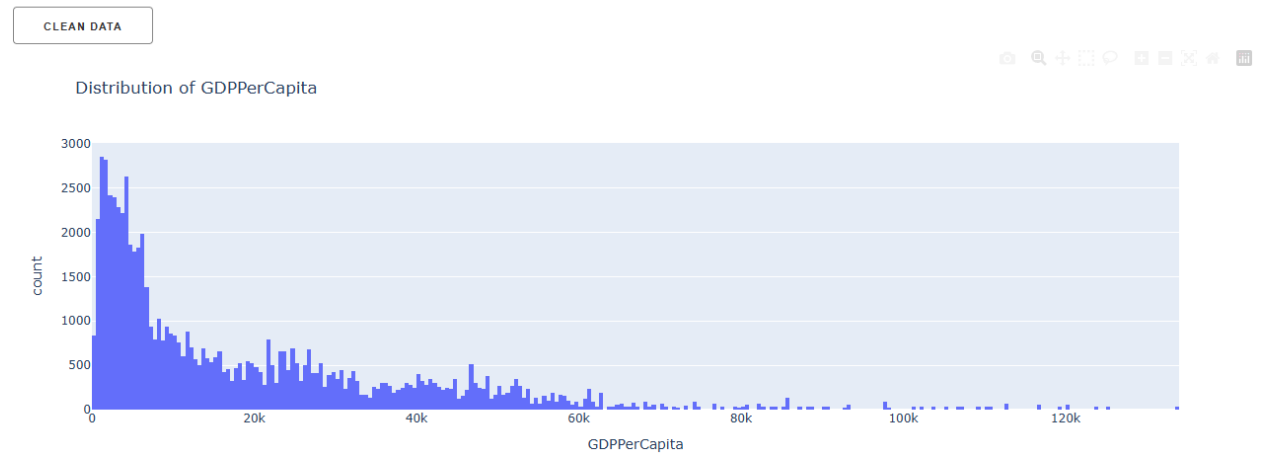


Figure 2.

Outlier detection and removal.

Outlier Detection and Removal Using IQR

Outliers have been removed. Data points outside the range [-79.00, 137.00] were removed.



Figure 3.

PCA

Dimensionality Reduction with PCA

9 components are required to retain 95% of the variance.

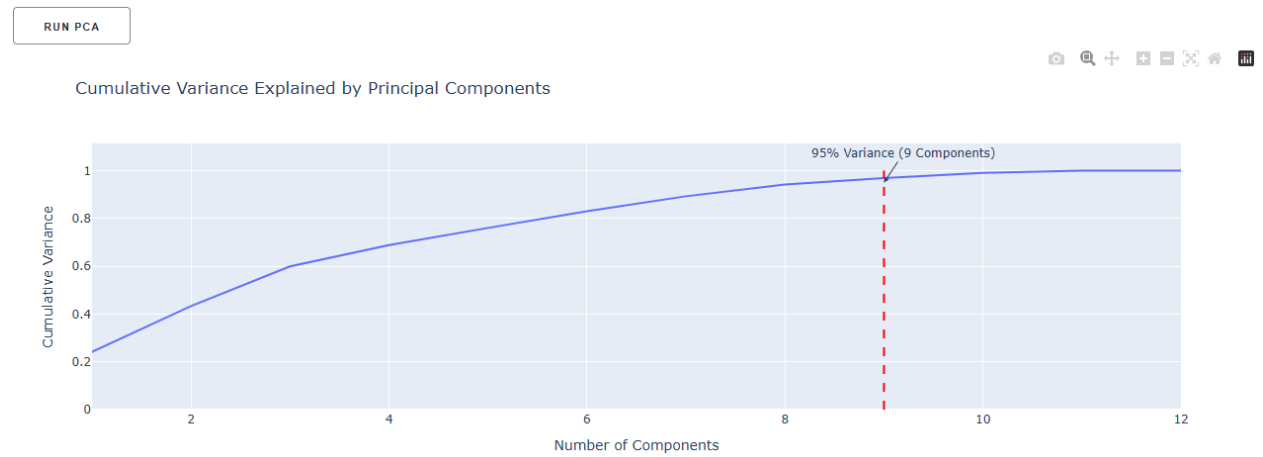


Figure 4.

Normality test.

Normality Test and Transformation

The column Population is not normally distributed (p-value = 0.0000). A log transformation has been applied.

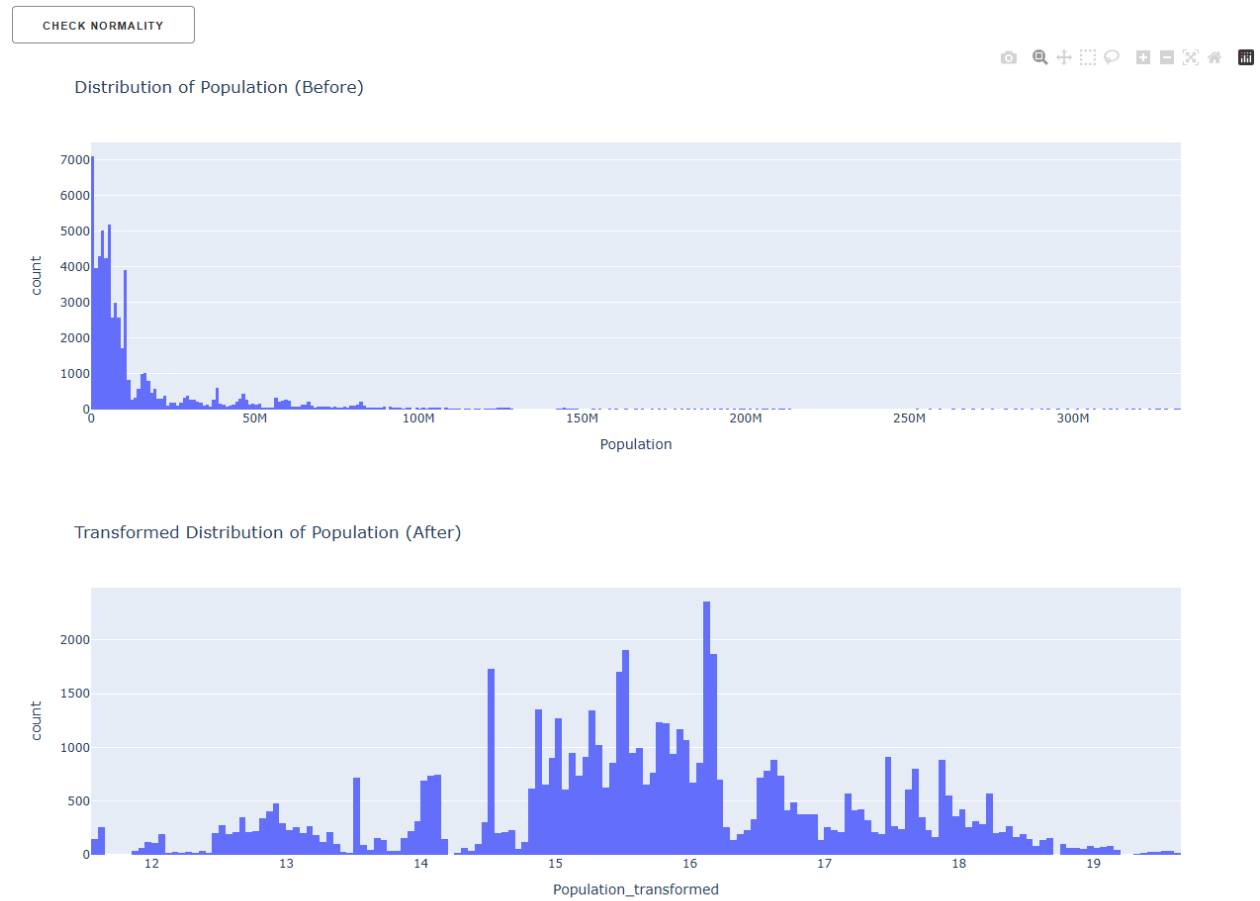


Figure 5.

Log and min max transformation.

Apply Data Transformation

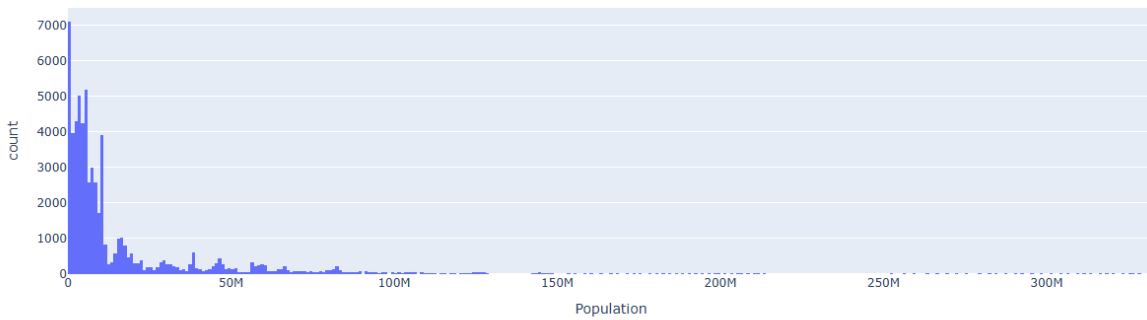
Select Transformation Type:

Log Transformation ✕ ▾

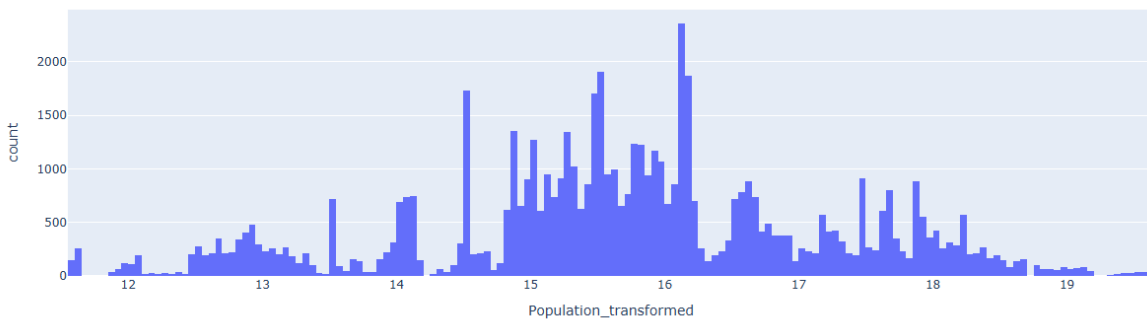
Log transformation applied. Transformation applied to Population.

APPLY TRANSFORMATION

Distribution of Population (Before)



Distribution of Population (After Log)



Apply Data Transformation

Select Transformation Type:

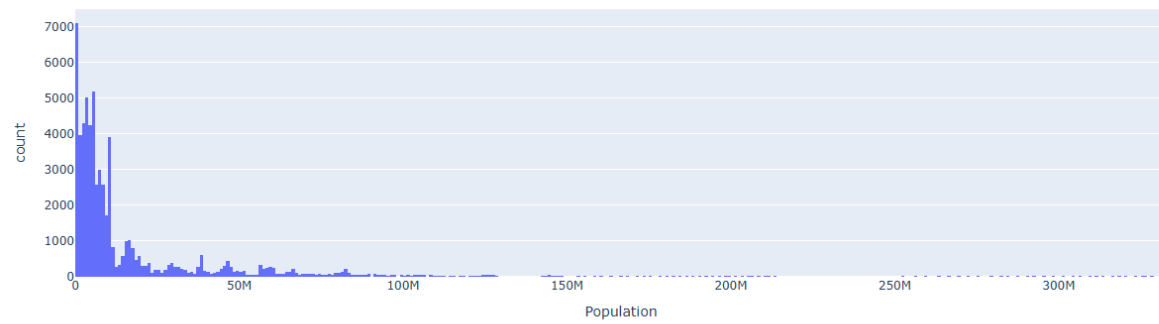
Min-Max Scaling

✕

Min-max scaling applied. Transformation applied to Population.

APPLY TRANSFORMATION

Distribution of Population (Before)



Distribution of Population (After Minmax)

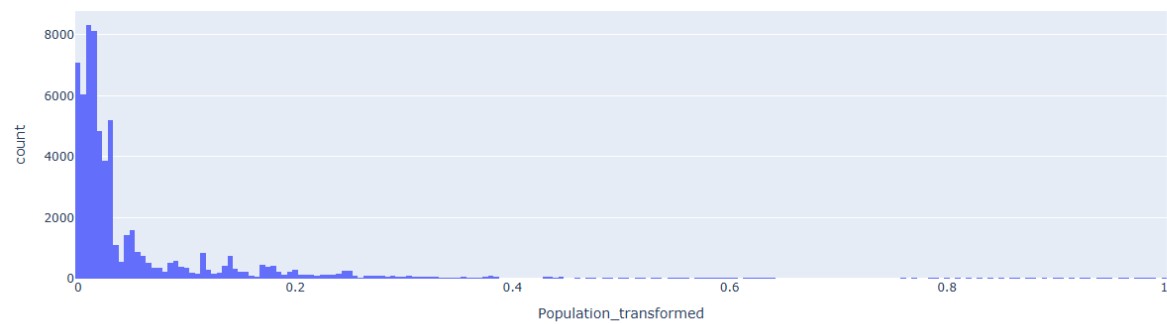


Figure 6.

Statistics.

Dataset Statistics

Select Feature for Statistics:

SuicideCount

✕ ▼

Mean: 22.11

Median: 9.00

Standard Deviation: 29.18

Minimum: 0.00

Maximum: 137.00

25th Percentile: 2.00

75th Percentile: 32.00

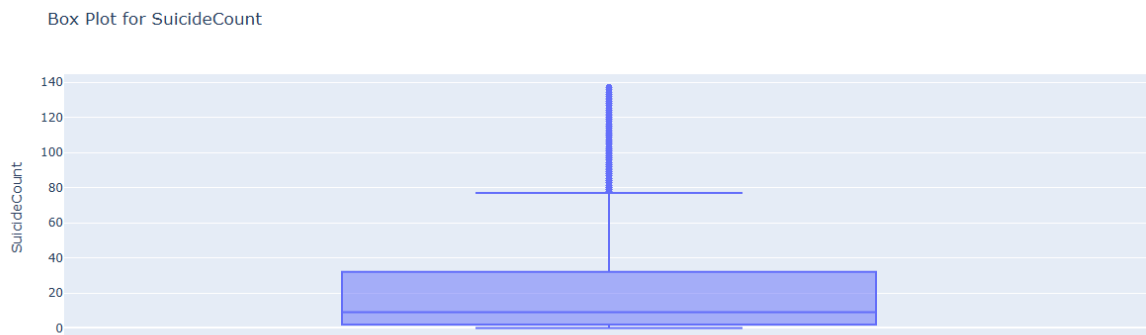


Figure 7.

Line plot.

Dynamic Plotting for All Features

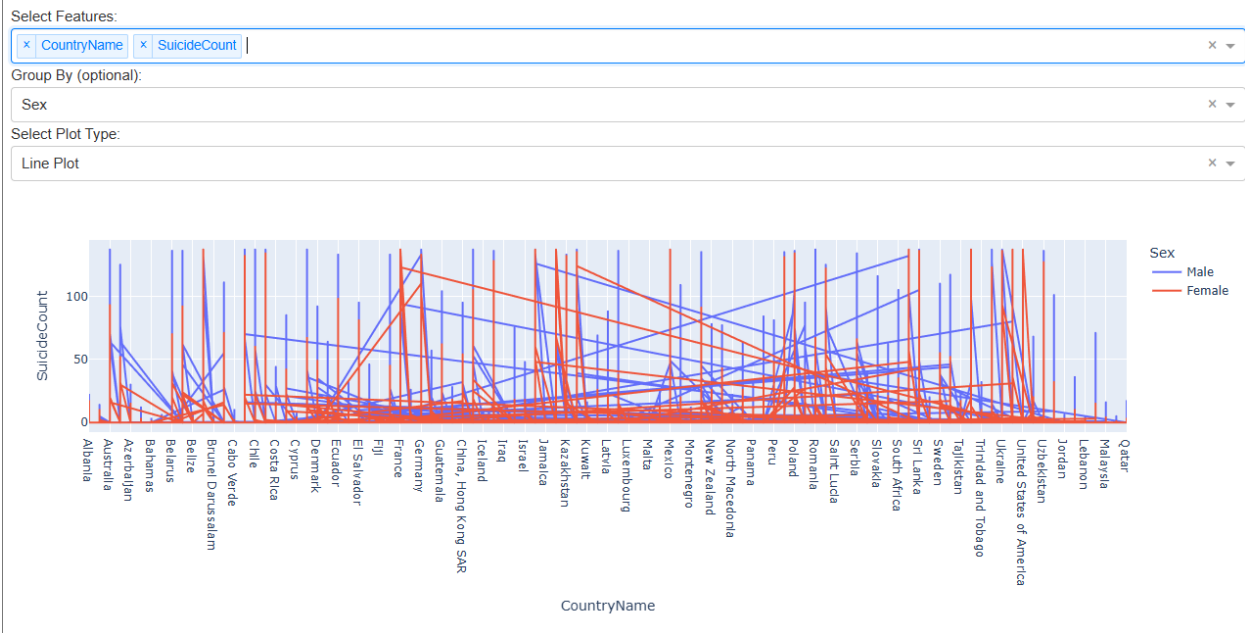


Figure 8.

Grouped bar plot.

Dynamic Plotting for All Features

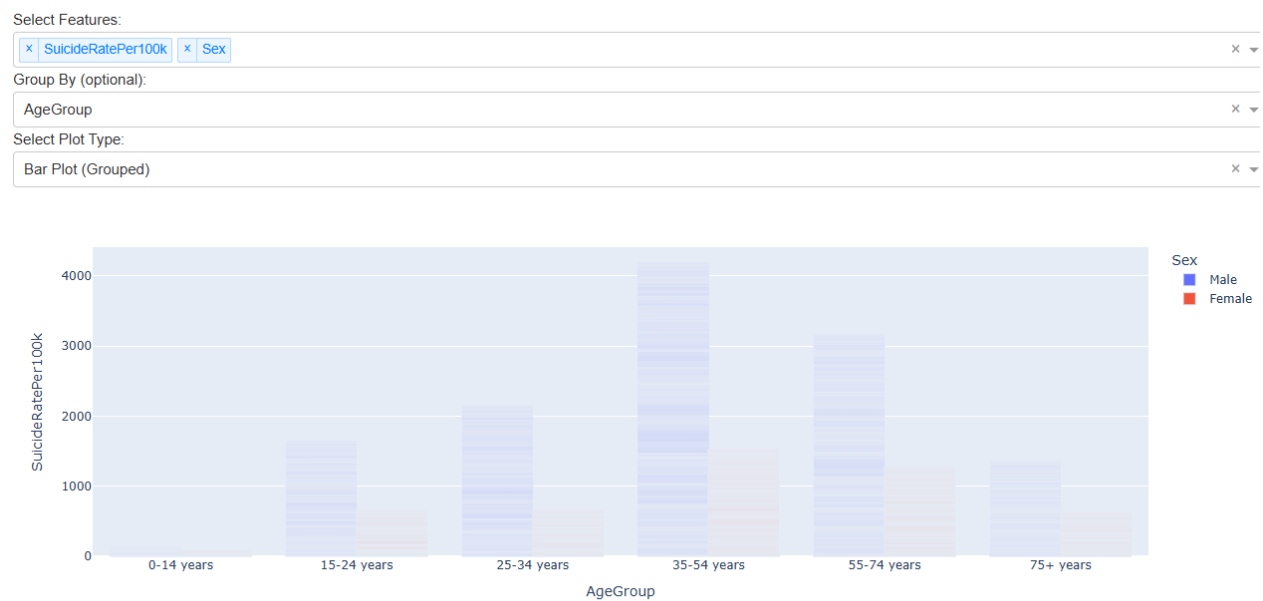


Figure 9.

Count plot.



Figure 10.

Pie chart.

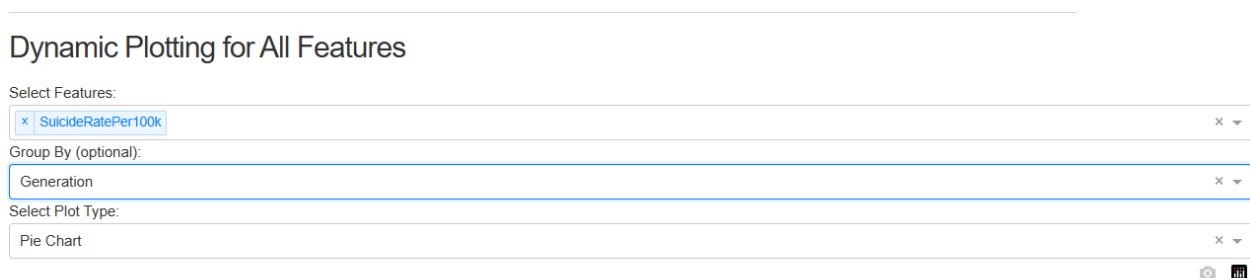


Figure 11.

Dist. Plot.

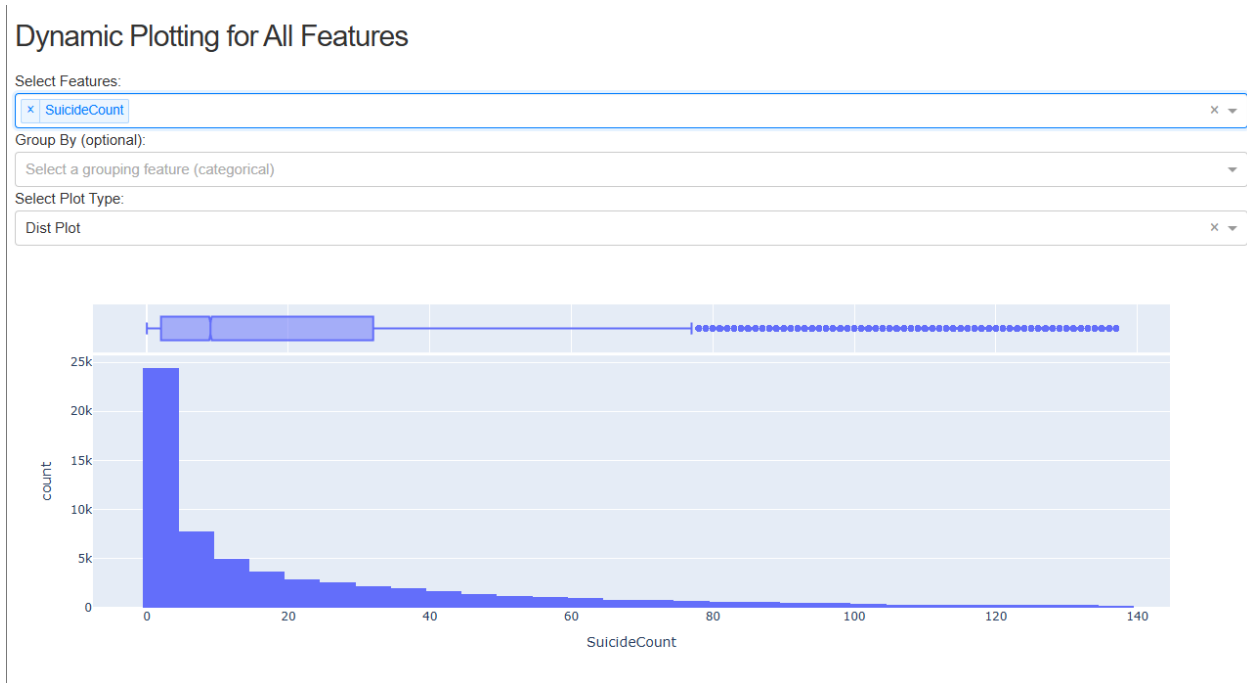


Figure 12.

Pair plot.

Dynamic Plotting for All Features



Figure 13.

Heat map.

Dynamic Plotting for All Features

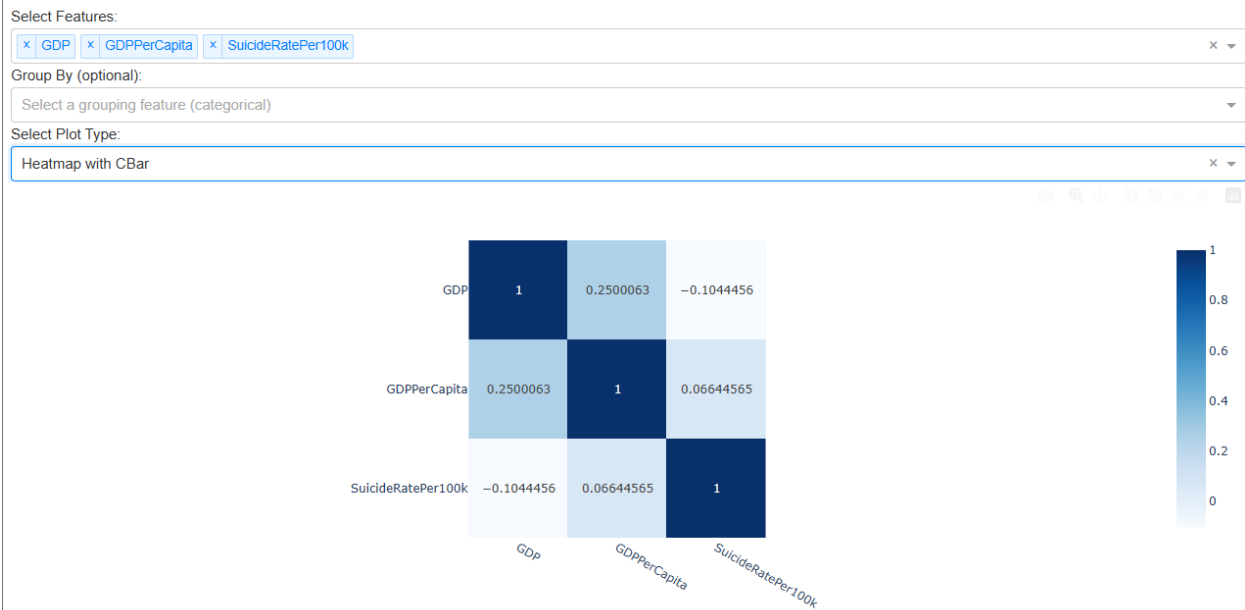


Figure 14.

Qq plot.

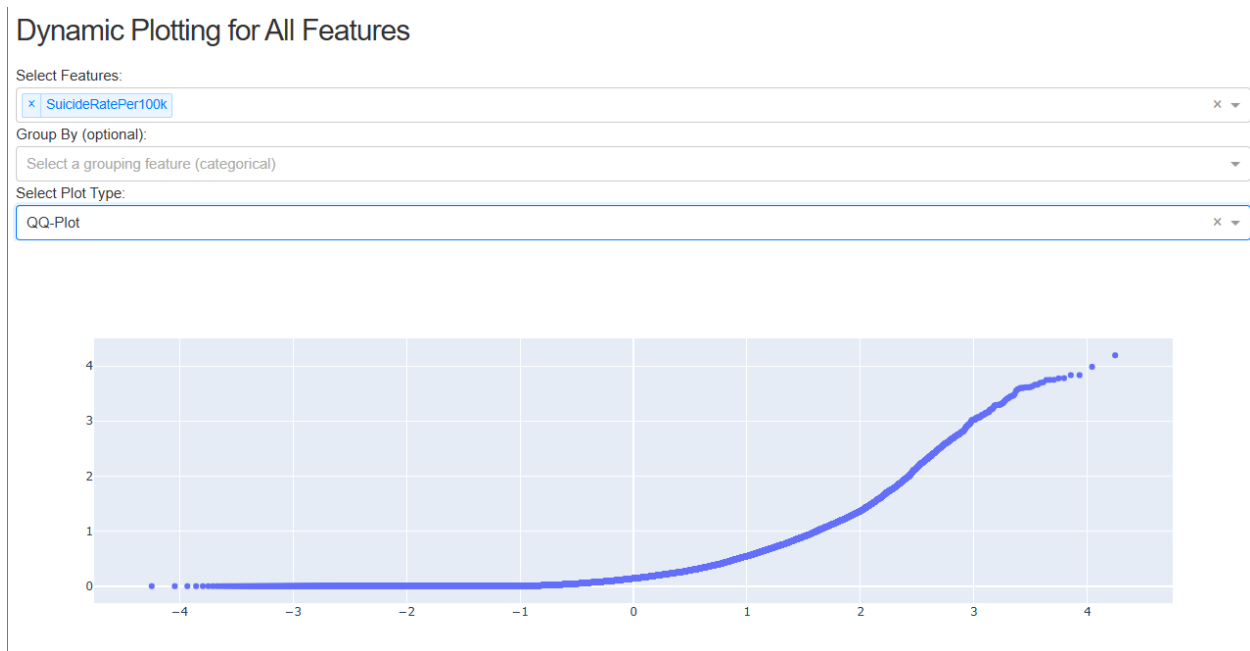


Figure 15.

KDE plot.

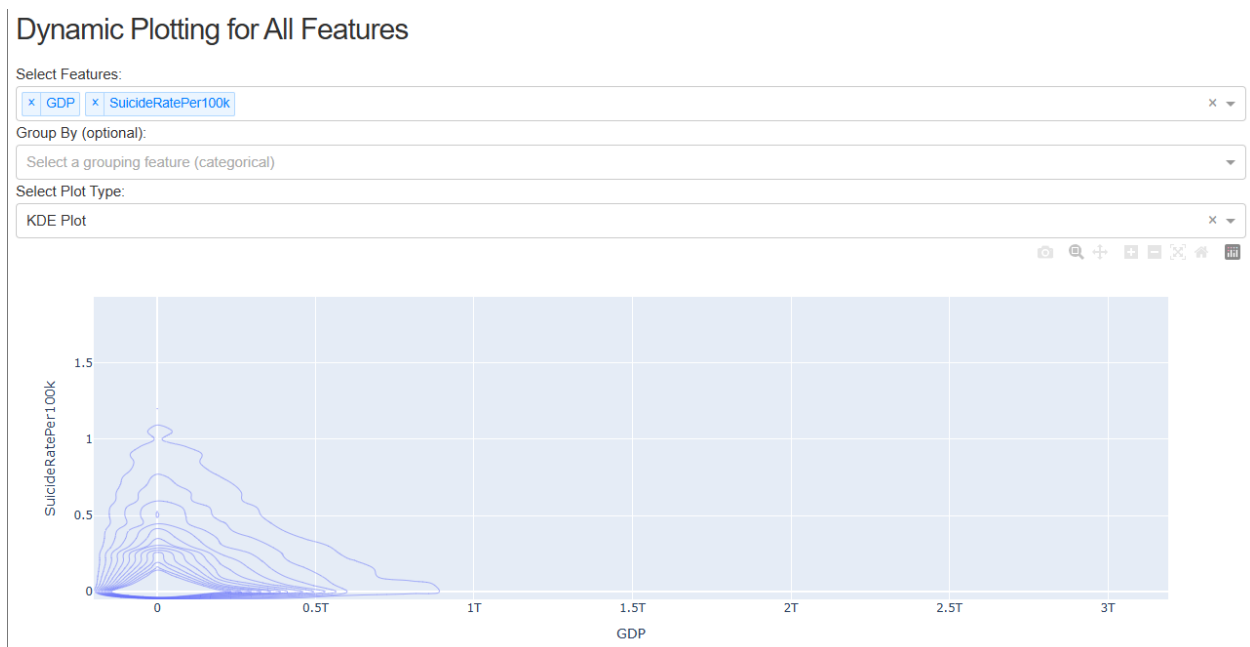


Figure 16.

Box plot.

Dynamic Plotting for All Features

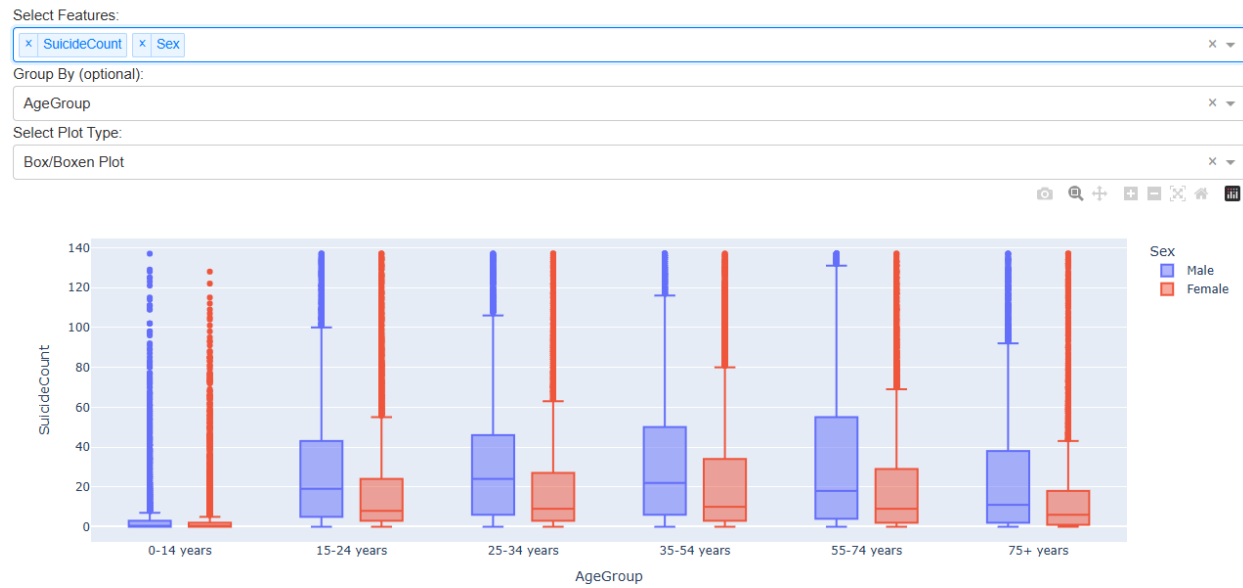


Figure 17.

Area plot.

Dynamic Plotting for All Features

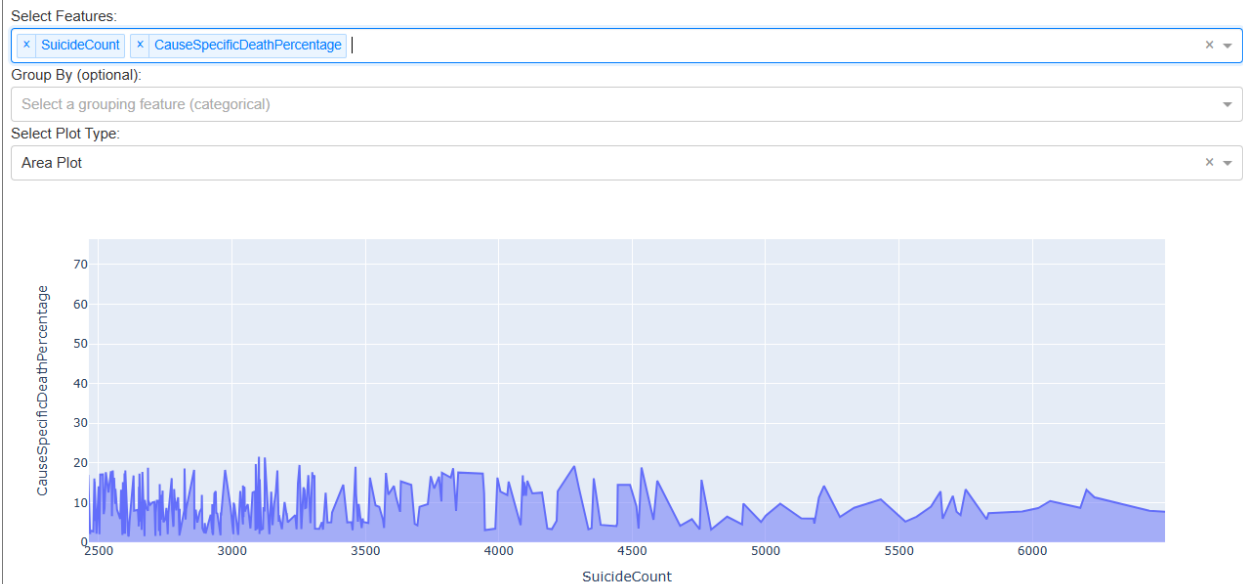


Figure 18.

Violin plot.

Dynamic Plotting for All Features

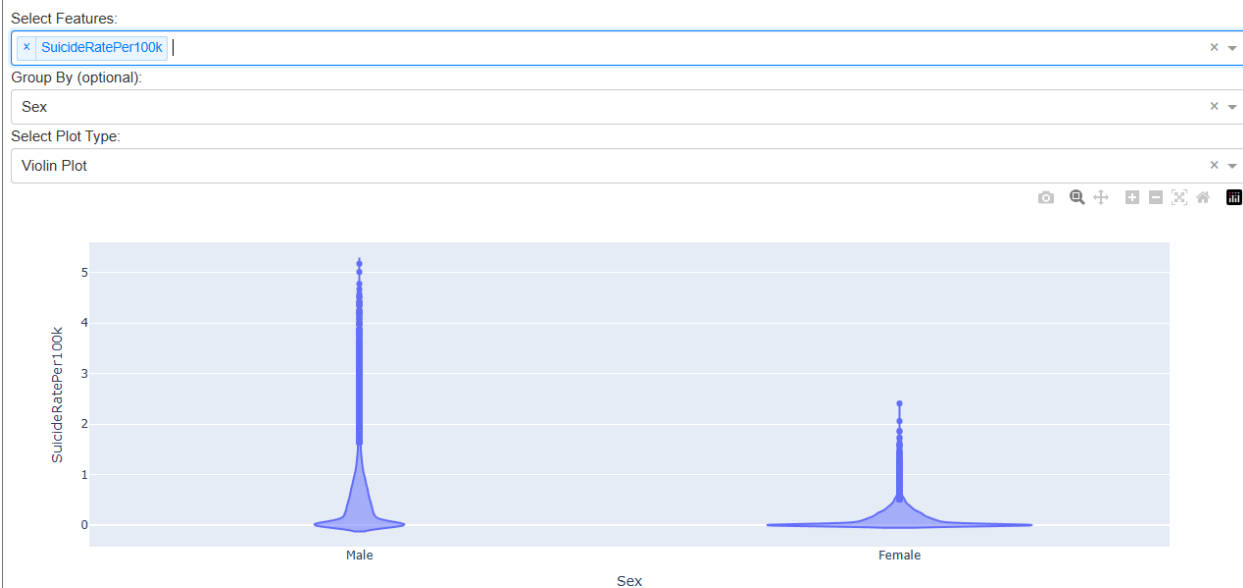


Figure 19.

Rug plot.

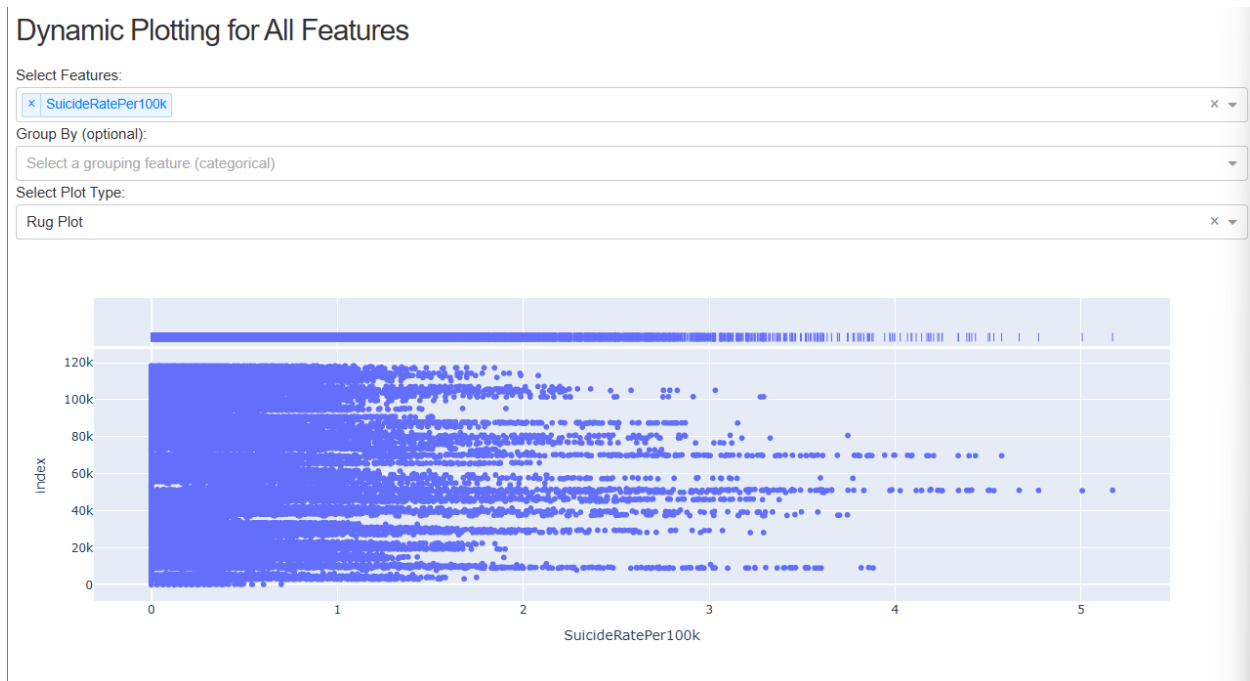


Figure 20.

3D plot.

Dynamic Plotting for All Features

Select Features:

x GDP

x InflationRate

x SuicideRatePer100k

x

Group By (optional):

Select a grouping feature (categorical)

Select Plot Type:

3D Plot

x



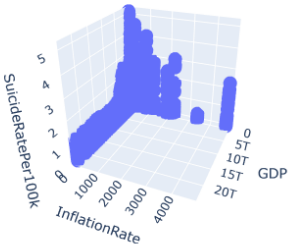


Figure 21.
Contour plot.

Dynamic Plotting for All Features

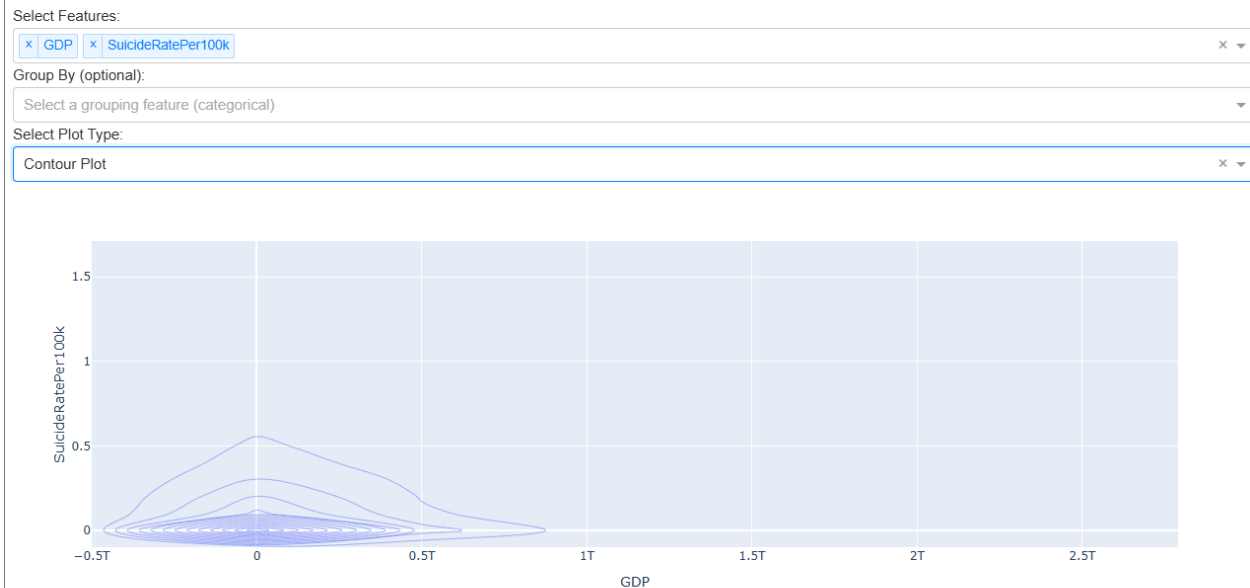


Figure 22.

Hexbin plot.

Dynamic Plotting for All Features

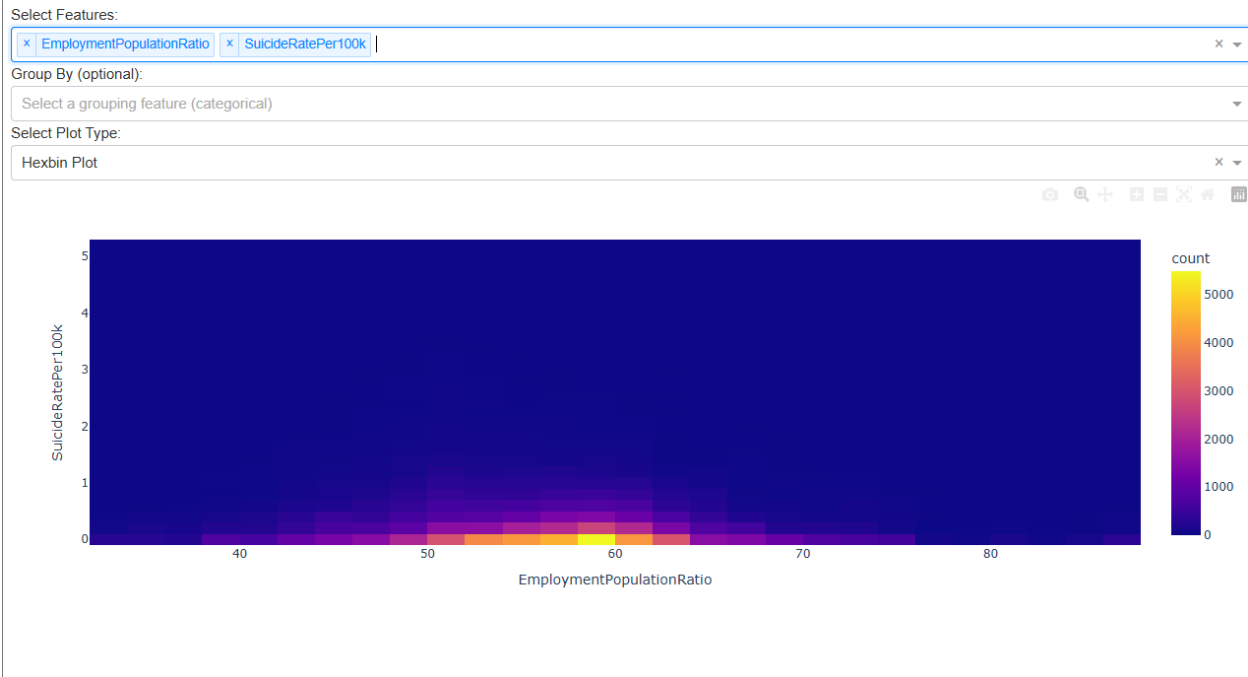


Figure 23.

Strip plot.

Dynamic Plotting for All Features

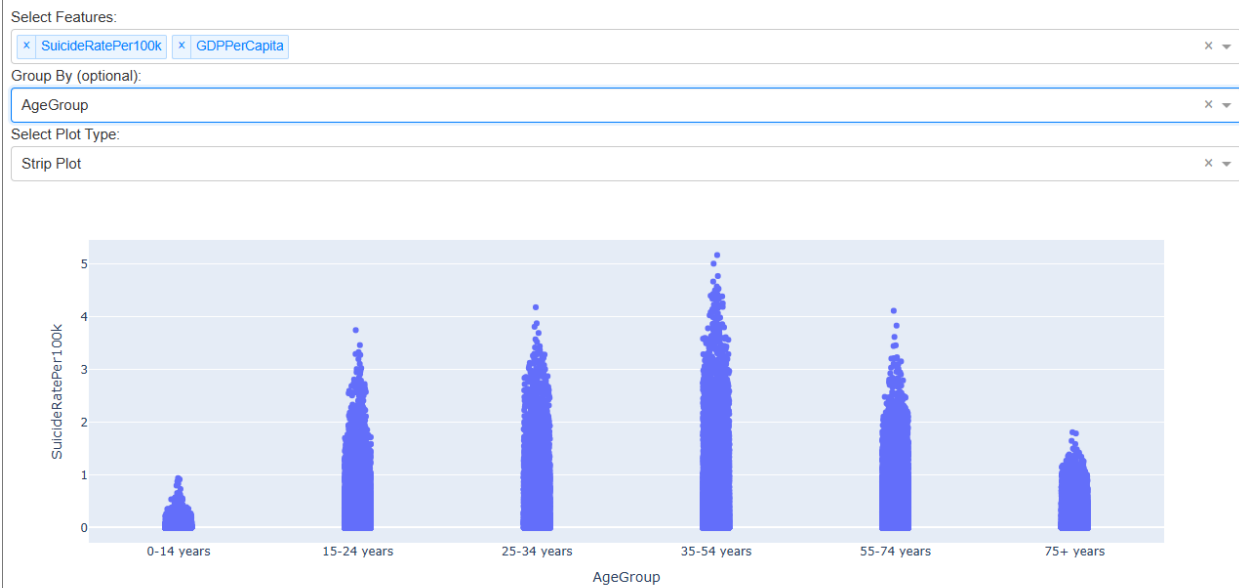


Figure 24.

Subplot.



Abstract

Suicide is a serious topic that has been talked about throughout the history of time. This paper explores the world of data visualization using python coding language. It demonstrates the process of visualizing selected dataset in various visualization techniques. Applying methods such as data cleaning, outlier detection and deletion, feature selection using PCA, statistical analysis, and finally displaying all that information on a user-friendly dashboard.

Keywords: information visualization, outlier detection, PCA, statistics.

Introduction

This study shows the process of understanding and visualizing a selected data set using python coding language. This study goes through various techniques, including exploratory data analysis, data processing, statistical analysis and visualizing data using an interactive dashboard made by Dash python library.

Description of the Dataset

The dataset contains 118,560 rows of observation and 18 columns of features. The dependent variable of this data set is SuicideCountPer100k. The independent variables include GDP, population, socioeconomic, employment rate and inflation rates. This data set is relevant in industry because of its importance in understanding the trends in public health and socioeconomic factors. The independent variable often time do play a significant role in prediction the trends in real world application, like mental health interventions.

Data pre-processing

This dataset was nowhere near ready for further analysis, the pre-processing contains using various techniques included in Python library Pandas. Non existing elements and duplicated observations were deleted (*figure 1*). The plot shows one of the features after the dataset was cleaned.

Outlier detection and removal

In this phase, outliers were detected using the IQR method and visualized with a box plot (*figure 2*). The after removal of outlier plot shows that values for Suicide Count beyond the dedicated IQR range were removed to better the quality of the dataset for future processing and analyzing.

Principal Component Analysis

After removing the outliers in previous phase. PCA was conducted, the functionality of PCA was to reduce feature dimensions while retaining 95% of the variance. 9 out of 12 features suffice the variance threshold, and 3 did not (*figure 3*). The conditional number and singular values were within the acceptable range.

Normality Test

The normality test was done by Shapiro-Wilk test. The feature Population was not gaussian. This resulted in a log transformation of the feature (*figure 4*).

Data Transformation

Data transformation was done by log transformation and min max transformation on the feature Population (*figure 5*). After transformation, the normality was improved, and this enhanced the dataset's ability to further analysis.

Heatmap

The heatmap (*figure 13*), highlighted the relationship between three variables, GDP, GDP per capita and suicide rate per 100k. Revealing minor correlation between suicide count and GDP per capita.

Statistics

The statistical portion of the dashboard (*figure 6*) includes, mean, median, and standard deviation. They were calculated for key features selected from the dropdown menu.

Data Visualization

The visualizations were listed in this paper through figure 7 to figure 23. Each observation will be covered in the **Observation** section.

Subplot

The subplot in *figure 23* shows the global picture of suicide rates. From the subplot, Europe has higher rates compared to other regions (*plot 1 in figure 23*). Suicide rates have fluctuated across years, but still with noticeable peaks in certain regions (*plot 2 in figure 23*). The record of suicide was mostly distributed in the Europe region (*plot 3 in figure 23*). The final plot (*plot 4 in figure 23*) shows that male has much higher suicide records than female counterparts.

Dashboard



Data cleaning

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	-------------------------------	--------------------------	-----------------	---------------------	------------	------------------

Outlier Detection and Removal Using IQR

Outliers have been removed. Data points outside the range [-79.00, 137.00] were removed.

REMOVE OUTLIERS



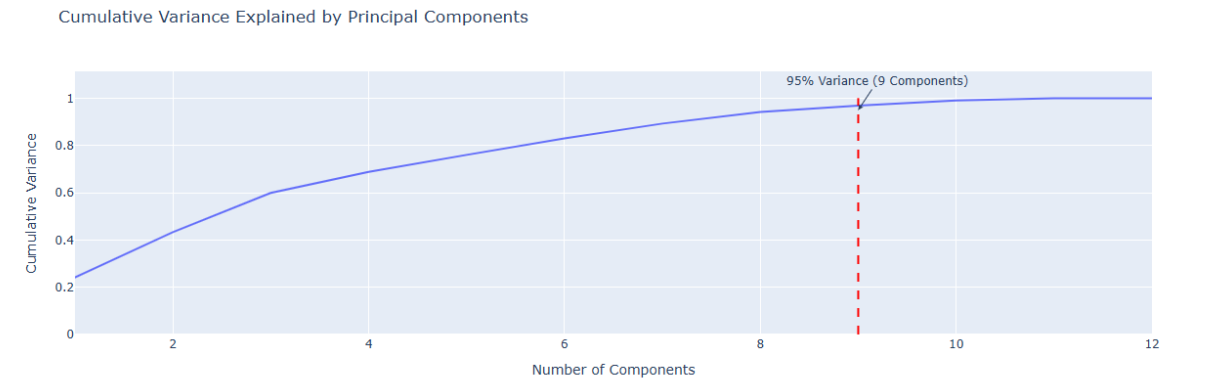
Outlier detection and removal

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

Dimensionality Reduction with PCA

9 components are required to retain 95% of the variance.

RUN PCA



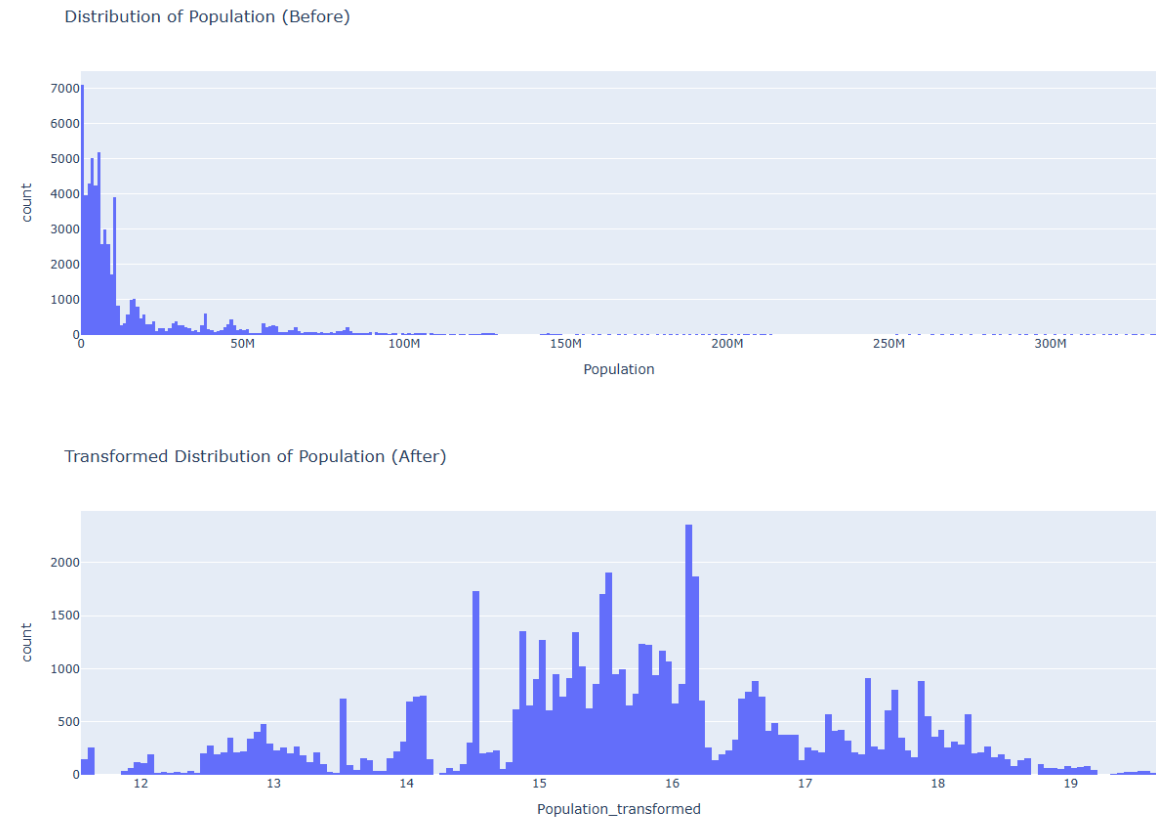
PCA

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

Normality Test and Transformation

The column Population is not normally distributed (p-value = 0.0000). A log transformation has been applied.

CHECK NORMALITY



Normality test

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

Apply Data Transformation

Select Transformation Type:

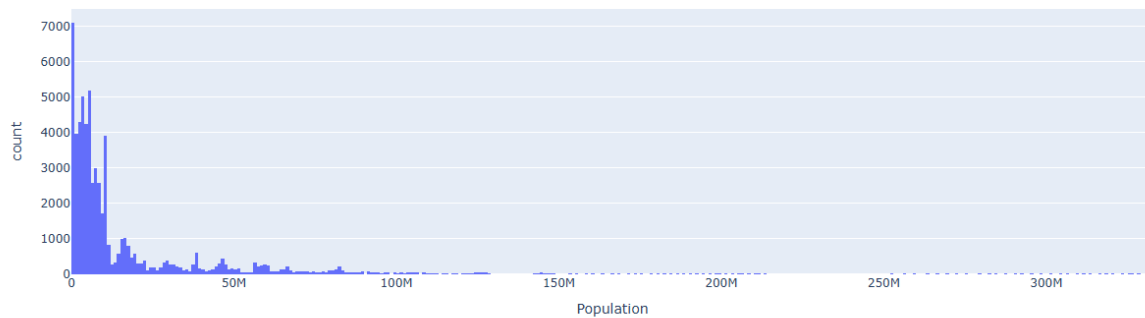
Log Transformation

X

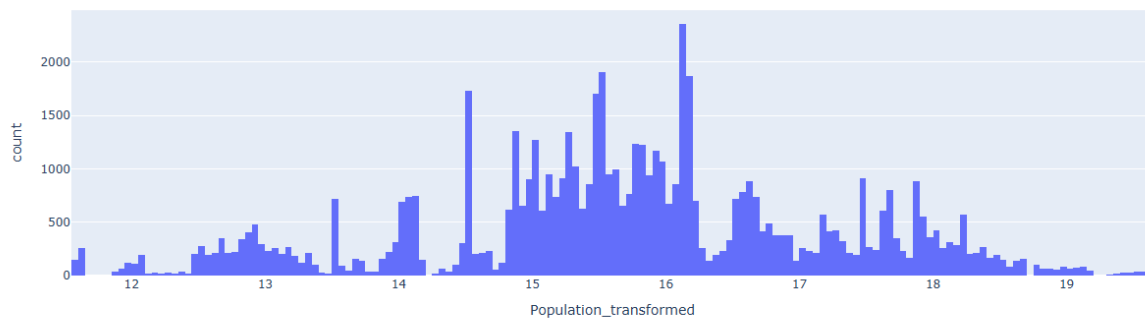
Log transformation applied. Transformation applied to Population.

APPLY TRANSFORMATION

Distribution of Population (Before)



Distribution of Population (After Log)



Data transformation

Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

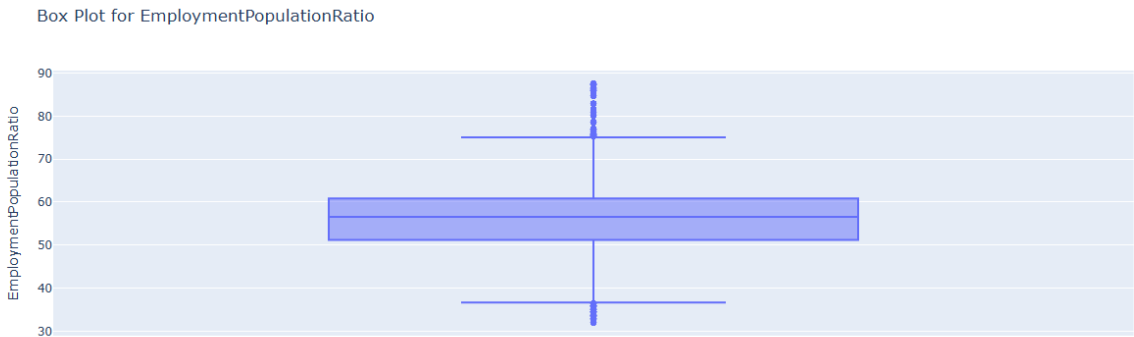
Dataset Statistics

Select Feature for Statistics:

EmploymentPopulationRatio

x

Mean: 56.15
Median: 56.60
Standard Deviation: 7.96
Minimum: 32.03
Maximum: 87.52
25th Percentile: 51.19
75th Percentile: 60.87



Statistic

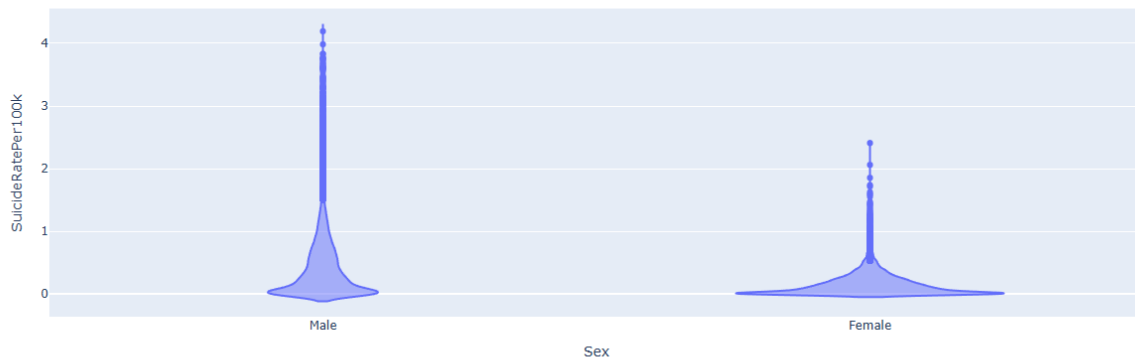
Data Cleaning	Outlier Detection and Removal	Dimensionality Reduction	Normality Tests	Data Transformation	Statistics	Dynamic Plotting
---------------	----------------------------------	-----------------------------	-----------------	------------------------	------------	------------------

Dynamic Plotting for All Features

Select Features:

Group By (optional):

Select Plot Type:



Dynamic plotting

Observation

Line plot (*figure 7*), this plot shows major variation in suicide count across different countries, there is a clear distinction between suicide records for males and females, with male having higher records. Some countries show significant spike meaning there may be some cultural aspects for that to happen.

Grouped bar plot (*figure 8*), this plot shows that suicide rate per 100k population increases significantly in older age groups, the peak is in the 35-54 years age group category. Male suicide still is significantly higher than female record counts in all age group.

Count plot (*figure 9*), this plot shows Generation X and Baby boomers have the highest count of suicide rate, and generation like Silent generation and generation Alpha have the lower suicide rate records.

Pie chart (*figure 10*), this plot shows the generation X and baby boomers combined for most suicide records, at 33 percent and 25 percent.

Dist. Plot (*figure 11*), this plot shows the distribution of suicide count is highly skewed to the left, with most data points at lower values, this might indicate the need for data transformation.

Pair plot (*figure 12*), this plot shows no strong liner correlation with suicide count and GDP, GDP per capita, suggesting that suicide rates might depend on other features.

Heat map (*figure 13*), this plot indicates a weak correlation between suicide rates per 100k and GDP per capita, suggesting that suicide rates might depend on other features.

QQ plot (*figure 14*), this plot shows deviation from the diagonal line, indicating that suicide rate per 100k does not follow a normal distribution.

KDE plot (*figure 15*), this plot shows clustering of data point at lower GDP and suicide rate values, the density decreases as GDP increases, suggesting the wealthier country might have lower suicide rate records.

Box plot (*figure 16*), this plot shows male suicide counts are consistently higher across all the age groups, reinforcing earlier observation.

Area plot (*figure 17*), this plot shows a decreasing trend of cause specific percentages is visible as suicide count increases, indicating that higher suicide counts are associated with a bigger range of contributing factors.

Violin plot (*figure 18*), this plot shows males exhibit a wider range and higher density of suicide rates compared to females.

Rug plot (*figure 19*), this plot shows a dense clustering of data points at lower suicide rates, some high outliers are visible, meaning there might be some extreme cases.

3D plot (*figure 20*), this plot shows no linear relationship between suicide rates, GDP, and inflation rates, but higher GDP values are associated with lower suicide rates.

strip plot (*figure 23*), this plot shows with the 35-54 years age group having the widest spread, and higher GDP per capita values are associated with lower suicide rates.

Conclusion

From this study, all the graphs show important insights. The correlation between GDP and suicide rates are not as clear, but the graphs suggested that lower GDP valued countries have higher suicide rates compared to countries with higher GDP values. The dashboard designed with python was user friendly, but sometimes it takes a while to render the graph, suggesting a need for higher computational power. In conclusion, suicide is a serious topic, no amount of data observed can help to reduce the cases of suicide, make sure you are being vigilant and help anyone that is close to you if they are having a hard time.

Appendix

The python files, one for phase I, and another for phase II will be submitted on canvas separate from this paper.

References

Python Libraries: *Pandas, NumPy, plotty and other various ML libraries.*

Various passed information visualization assignments and in-class coding samples.

Dataset Source:

“Suicide Rates & Socioeconomic Factors (1990 - 22).” Onyango, Kaggle, 1 Mar. 2024,

<https://www.kaggle.com/datasets/ronaldonyango/global-suicide-rates-1990-to-2022/>.