# Data Analytics

## PYTHON

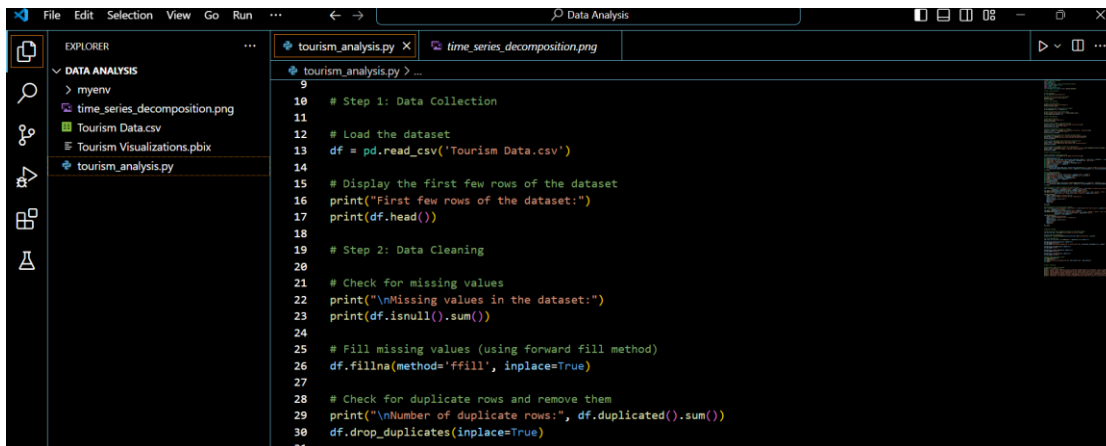Reem Hamraz | Mindshare Yuva | 29.06.24

# Introduction

Hello!!

I have created this document to elaborate upon and present my analysis of a hypothetical tourism dataset that spans from 2015 to 2020. This dataset includes information on the number of tourists, revenue, and average spending across various countries. I've used Python for data analysis and visualization, as well as Power BI to enhance the presentation with interactive visuals.

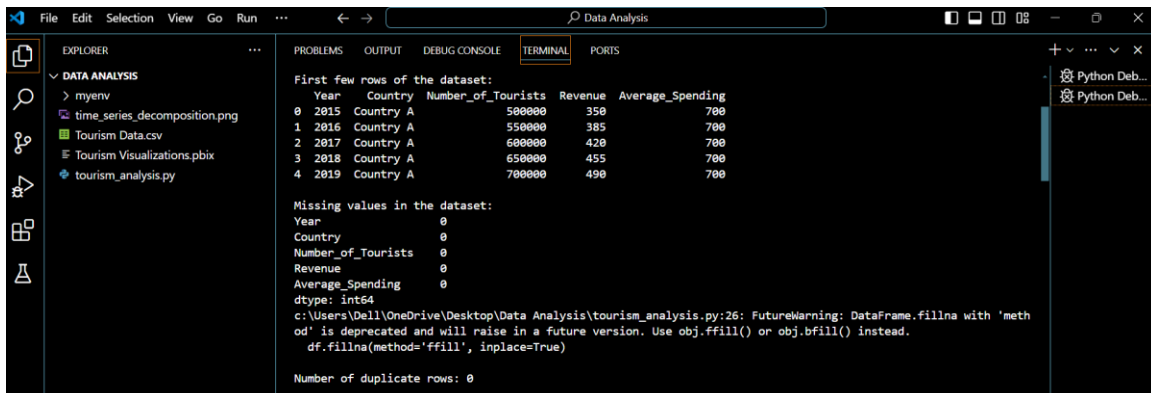## 1. DATA ANALYSIS AND VISUALIZATION IN PYTHON

1. **Data Collection and Cleaning:** First, I imported the dataset using the Pandas library and performed initial exploration by displaying the first few rows. This helped me understand the structure and contents of the data. After that, I checked for and handled any missing values and duplicate rows to ensure data quality.
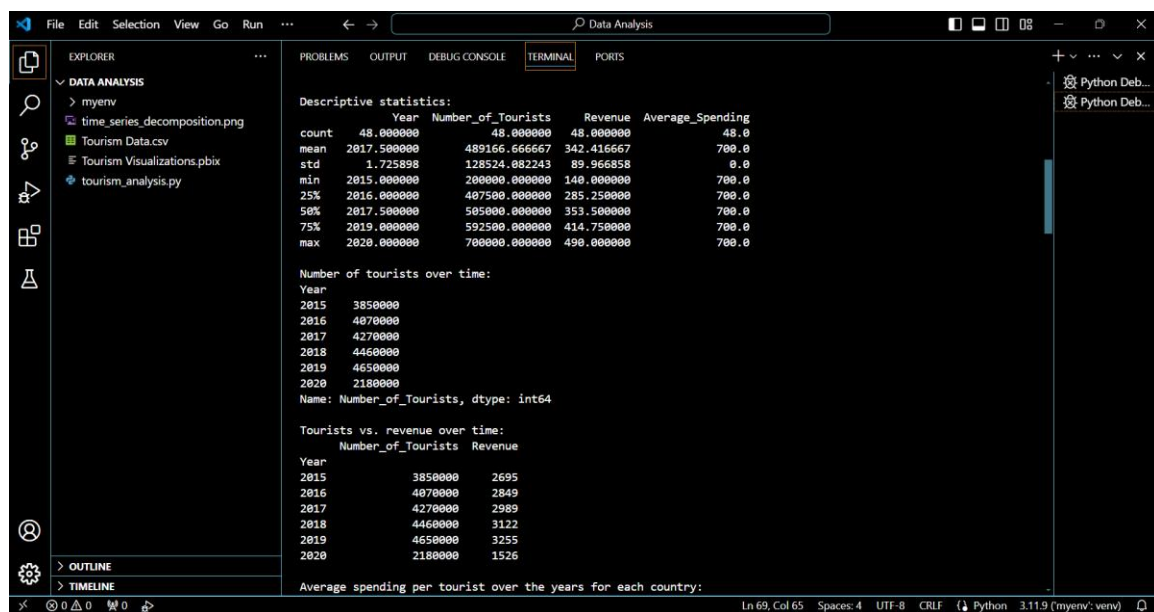
*Code-*



*Output-*

2. **Data Analysis:** I conducted a descriptive statistical analysis to summarize the data. This included calculating the total number of tourists, total revenue, and average spending per tourist over the years for each country. Additionally, I analyzed trends such as the yearly change in the number of tourists for each country and the relationship between the number of tourists and revenue.

*Code-*

```python
32    # Step 3: Data Analysis
33
34    # Descriptive statistics
35    print("\nDescriptive statistics:")
36    print(df.describe())
37
38    # Analyze trends: number of tourists over time
39    tourists_over_time = df.groupby('Year')['Number_of_Tourists'].sum()
40    print("\nNumber of tourists over time:")
41    print(tourists_over_time)
42
43    # Analyze relationships: tourists vs. revenue
44    tourists_vs_revenue = df.groupby('Year')[['Number_of_Tourists', 'Revenue']].sum()
45    print("\nTourists vs. revenue over time:")
46    print(tourists_vs_revenue)
47
48    # Average spending per tourist over the years for each country
49    avg_spending = df.groupby(['Year', 'Country'])['Average_Spending'].mean()
50    print("\nAverage spending per tourist over the years for each country:")
51    print(avg_spending)
52
53    # Total revenue generated by each country over the given period
54    total_revenue = df.groupby('Country')['Revenue'].sum()
55    print("\nTotal revenue generated by each country over the given period:")
56    print(total_revenue)
57
58    # Yearly change in the number of tourists for each country
59    yearly_change = df.groupby(['Country', 'Year'])['Number_of_Tourists'].sum().groupby(level=0).pct_change()
60    print("\nYearly change in the number of tourists for each country:")
```

*Output-*

```
Descriptive statistics:
              Year  Number_of_Tourists        Revenue  Average_Spending
count    48.000000           48.000000      48.000000              48.0
mean   2017.500000       489166.666667     342.416667             700.0
std       1.725898       128524.082243      89.966858               0.0
min    2015.000000       200000.000000     140.000000             700.0
25%    2016.000000       407500.000000     285.250000             700.0
50%    2017.500000       505000.000000     353.500000             700.0
75%    2019.000000       592500.000000     414.750000             700.0
max    2020.000000       700000.000000     490.000000             700.0

Number of tourists over time:
Year
2015    3850000
2016    4070000
2017    4270000
2018    4460000
2019    4650000
2020    2180000
Name: Number_of_Tourists, dtype: int64

Tourists vs. revenue over time:
      Number_of_Tourists  Revenue
Year
2015             3850000     2695
2016             4070000     2849
2017             4270000     2989
2018             4460000     3122
2019             4650000     3255
2020             2180000     1526

Average spending per tourist over the years for each country:
```
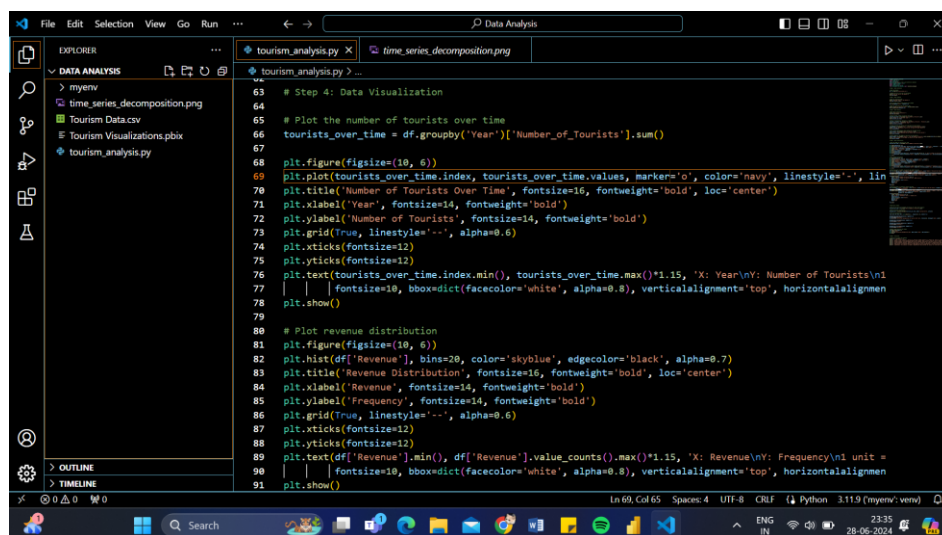
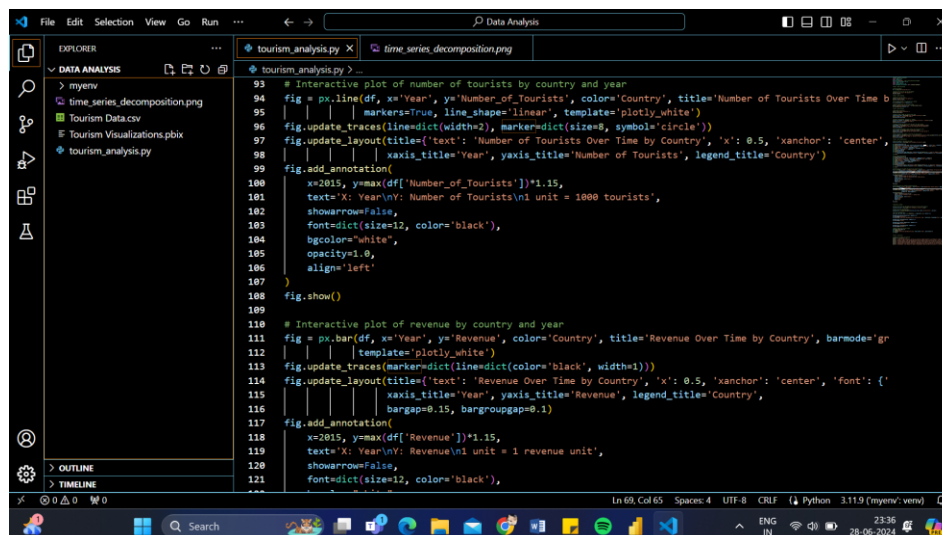3. **Data Visualization:** To visualize the data, I created several plots:

- **Line Chart:** This chart shows the number of tourists over time. It helps to visualize the trend and fluctuations in tourist numbers across the years.
- **Histogram:** The histogram displays the distribution of revenue, highlighting how revenue values are spread across the dataset.

Furthermore, I enhanced these visualizations by adding grid lines, borders, and customizing the aesthetics to make them more appealing and easier to interpret.

*Code-*

*Output-*

4. **Time Series Decomposition:** I performed a time series decomposition to break down the number of tourists into three components: trend, seasonality, and residuals. This decomposition helps in understanding the underlying patterns in the data. The trend component shows the long-term movement in the number of tourists, seasonality captures the periodic fluctuations, and residuals represent the random noise.
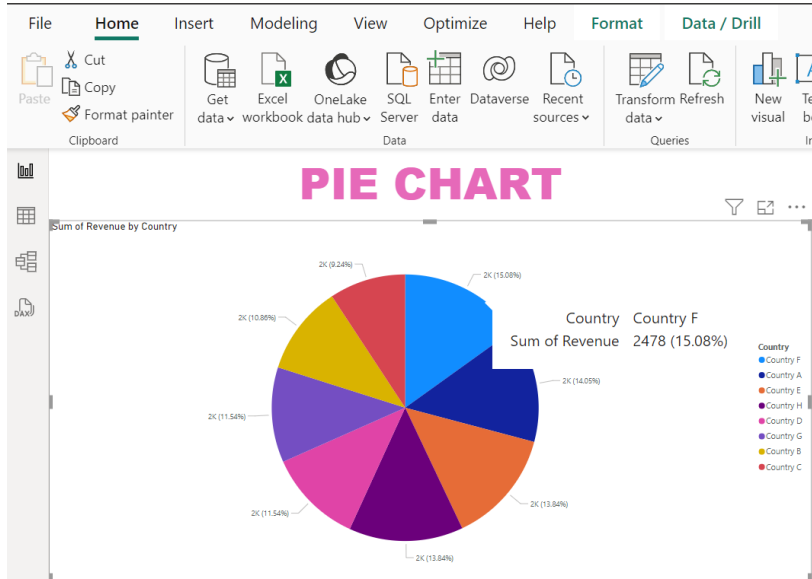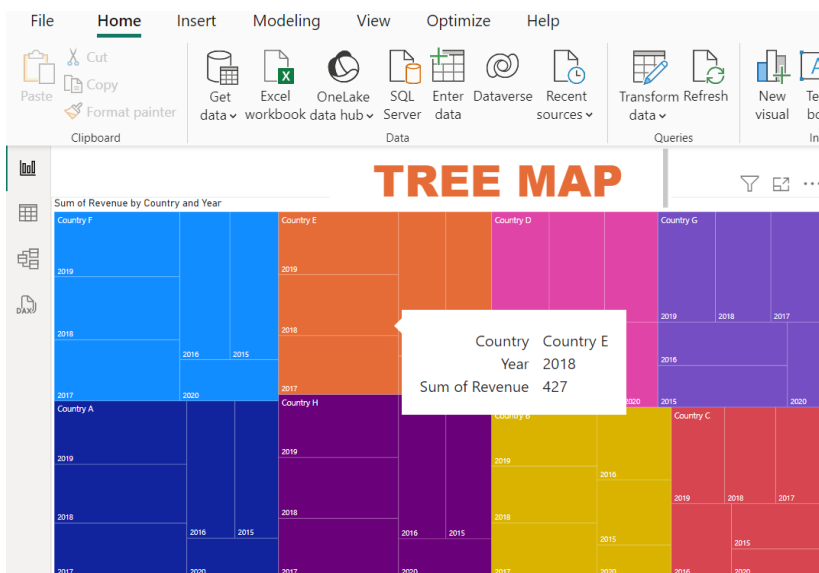
*Code-*



*Output-*

## 2. POWER BI VISUALIZATIONS

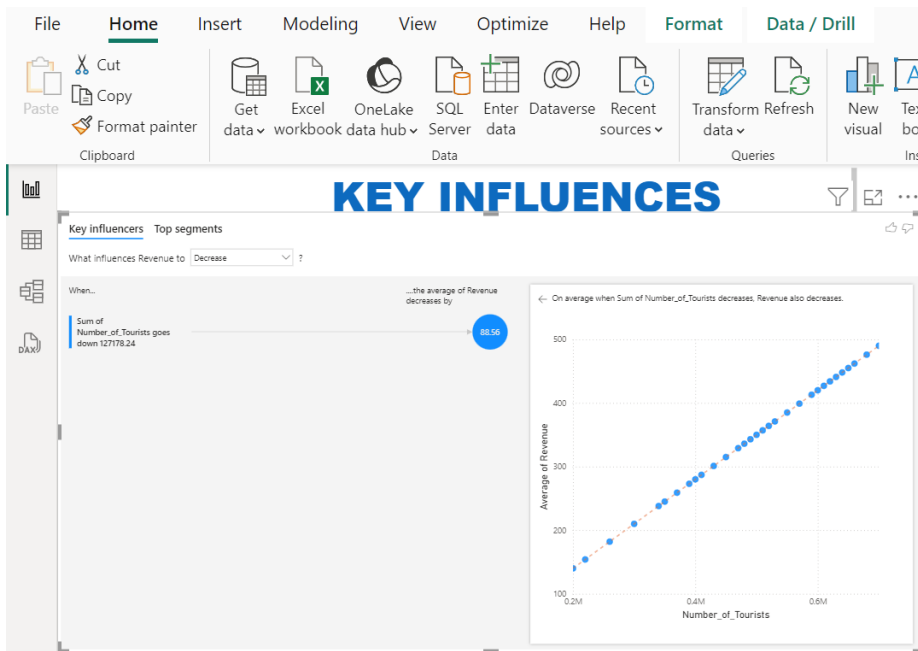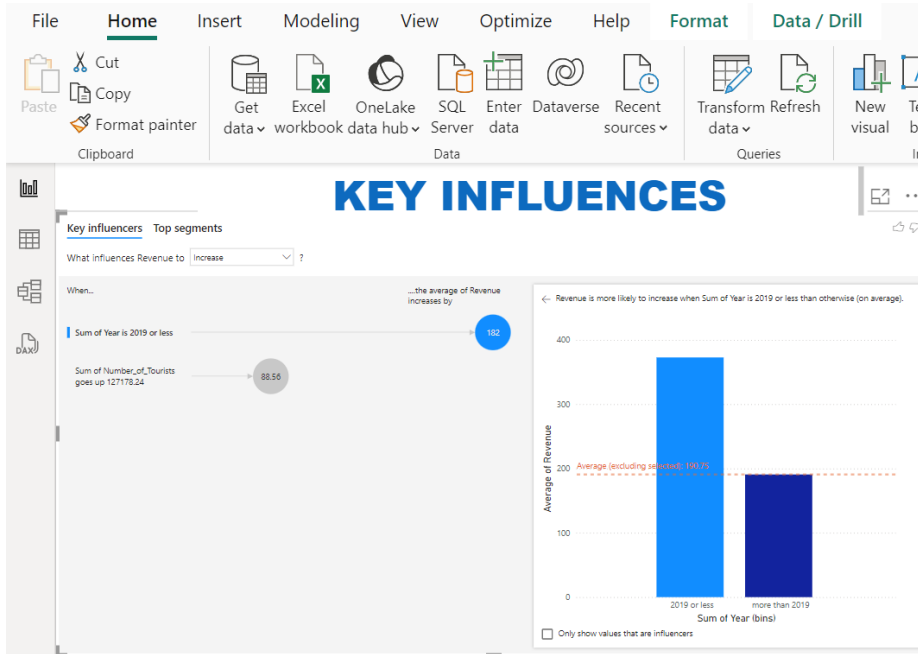Furthermore, I used Power BI to create interactive visualizations:

1. **Pie Chart:** The pie chart illustrates the revenue share by country. Each slice represents a country's contribution to the total revenue, making it easy to compare the relative performance of different countries.



2. **Treemap:** The treemap visualizes the contribution of each country to the total number of tourists. It uses a hierarchical layout to show the proportional size of each country's tourist numbers, providing a clear and immediate comparison.

3. **Key Influences:** The Key Influences visual helps to identify the factors that have the most significant impact on the number of tourists. It highlights which variables, such as revenue or specific years, influence tourist numbers the most. This visual is particularly useful for uncovering insights and driving data-driven decisions.

# Conclusion

Through this project, I gained valuable insights into the trends and patterns in the tourism world. The combination of Python for detailed analysis and Power BI for interactive visualizations has provided a comprehensive understanding of the data. I hope this document has demonstrated the power of combining these tools to analyze and visualize data effectively.



X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---X---