**A PROJECT REPORT**

On

# HOTEL BOOKING PREDICTION

*Submitted to CMREC (UGC Autonomous)*

*In Partial Fulfillment of the requirements for the Award of Degree of*

**BACHELOR OF TECHNOLOGY**

IN
## INFORMATION TECHNOLOGY

## MACHINE LEARNING

*Submitted By*

**MOHAMMED ARSH KHAN**                                    **228R1A1294**

*Under the guidance of*

**Mrs. M.JHANSI LAKSHMI**

Assistant Professor

Internal Guide

Mrs. M.JHANSI LAKSHMI

Assistant professor,

Department of IT

Head of the Department

Dr. MADHAVIPINGILI

Assoc Professor & HOD,

Department of IT

# **TABLE OF CONTENTS**

| **S. No.** | **Section** | **Page Number** |
|:---:|:---:|:---:|

# ABSTRACT

The hotel industry faces significant challenges due to booking cancellations and unpredictable customer behaviors. To address this, our project performs exploratory data analysis (EDA) on a comprehensive hotel booking dataset using Python libraries such as Pandas, NumPy, Matplotlib, and Seaborn. The dataset, which includes booking information for city and resort hotels, is cleaned and preprocessed to manage missing values and eliminate inconsistent data, such as bookings with no guests. The project replaces missing values intelligently using statistical techniques like mean imputation and mode substitution. Visualizations such as pie charts, bar plots, and heatmaps help uncover trends in booking status, lead time, and booking sources by country. These insights offer valuable indicators for predicting cancellation likelihood and optimizing hotel management strategies. Key findings include dominant countries in bookings, lead time distributions, and booking behaviors associated with cancellations. Though the project focuses primarily on EDA, it sets the stage for future predictive modeling using classification algorithms. By understanding customer behavior and booking trends, this analysis can support better forecasting, enhance revenue management, and reduce operational uncertainties. Ultimately, this project highlights the importance of data-driven decision-making in the hospitality industry through effective data preprocessing and visualization.

# **<u>INTRODUCTION</u>**

The hotel booking industry operates in a dynamic environment where customer behavior is influenced by various factors such as location, time of booking, and external economic conditions. Predicting and analyzing booking patterns is crucial for optimizing operations, reducing cancellation rates, and enhancing customer satisfaction. In this project, we focus on analyzing hotel booking data to identify trends and gain insights that can support future predictive analytics.The dataset used in this study contains extensive information about bookings made at both city hotels and resort hotels. We begin by loading and preprocessing the data using Python's Pandas library. Data cleaning involves handling missing values—replacing NaN entries in columns such as agent, company, and children with statistically derived values like zero or mean. Entries with no guests (adults, children, or babies) are removed as they do not represent valid bookings.The core of this project involves exploratory data analysis (EDA), a technique for summarizing the main characteristics of a dataset. Algorithms and techniques used include:Descriptive Statistics: For understanding distributions and central tendencies.Data Imputation: Using mean, mode, and conditional replacement for missing data.Data Visualization: Using Matplotlib and Seaborn to visualize patterns like lead time distribution, hotel type distribution, and top booking countries.Correlation Analysis: Using Pearson correlation matrix to evaluate relationships between numerical variables.This foundational analysis not only reveals booking trends but also provides a base for applying machine learning algorithms like logistic regression or decision trees in future predictive models.

## SOURCE CODE:

```
!pip install pycountry

Import numpy as np
Import pandas as pd
Import matplotlib.pyplot as plt
Import seaborn as sns
Import pycountry as pc

Pd.options.display.max_columns = None

## Importing Data
Data = pd.read_csv('hotel_bookings.csv')

#!pip install pycountry

Import numpy as np
Import pandas as pd
Import matplotlib.pyplot as plt
Import seaborn as sns
Import pycountry as pc

Pd.options.display.max_columns = None

## Importing Data
Data = pd.read_csv('hotel_bookings.csv')

## Show the first 5 rows of Data
Data.head()

## Copy the dataset
Df = data.copy()

## Find the missing value, show the total null values for each column and sort it in descending order
Df.isnull().sum().sort_values(ascending=False)[:10]

## Drop Rows where there is no adult, baby and child
Df = df.drop(df[(df.adults+df.babies+df.children)==0].index)

## If no id of agent or company is null, just replace it with 0
Df[['agent','company']] = df[['agent','company']].fillna(0.0)

## For the missing values in the country column, replace it with mode (value that appears most often)
Df['country'].fillna(data.country.mode().to_string(), inplace=True)
```

4

```
## For missing children value, replace it with rounded mean value
Df['children'].fillna(round(data.children.mean()), inplace=True)

# Distribution of hotel types
Hotel_counts = df['hotel'].value_counts()
Hotel_counts.plot.pie(autopct='%1.1f%%', figsize=(6, 6), title='Distribution of Hotel Types', ylabel='')
Plt.show()

# Top 10 countries with the most bookings
Top_countries = df['country'].value_counts().head(10)
Sns.barplot(x=top_countries.values, y=top_countries.index, palette='viridis')
Plt.title('Top 10 Countries with Most Bookings')
Plt.xlabel('Number of Bookings')
Plt.ylabel('Country')
Plt.show()

# Lead time distribution
Plt.figure(figsize=(8, 5))
Sns.histplot(df['lead_time'], kde=True, bins=50, color='blue')
Plt.title('Lead Time Distribution')
Plt.xlabel('Lead Time (days)')
Plt.ylabel('Frequency')
Plt.show()

# Count of reservation status
Sns.countplot(data=df, x='is_canceled', palette='cool')
Plt.title('Booking Status')
Plt.xlabel('Canceled (0=No, 1=Yes)')
Plt.ylabel('Count')
Plt.xticks([0, 1], ['Not Canceled', 'Canceled'])
Plt.show()
# Correlation heatmap (only numeric columns)
Plt.figure(figsize=(10, 8))

# Select only numeric columns
Numeric_df = df.select_dtypes(include=['number'])

# Compute the correlation matrix
Corr_matrix = numeric_df.corr()

# Plot the heatmap
Sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm')
Plt.title('Correlation Matrix')
Plt.show()
```
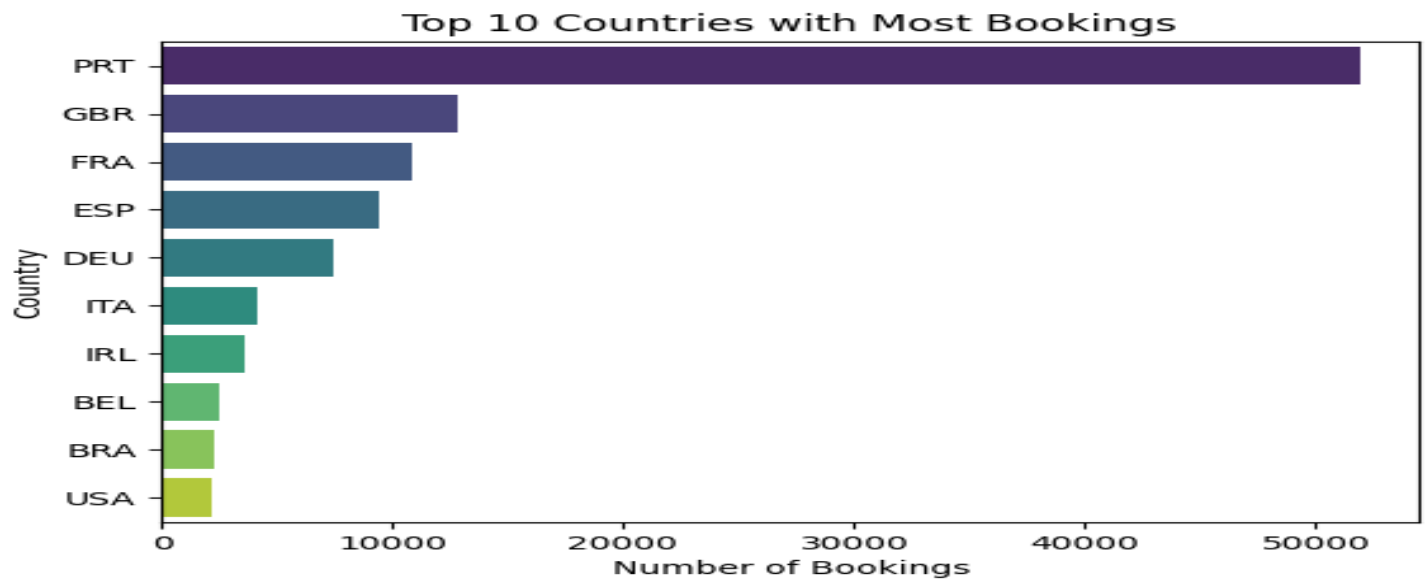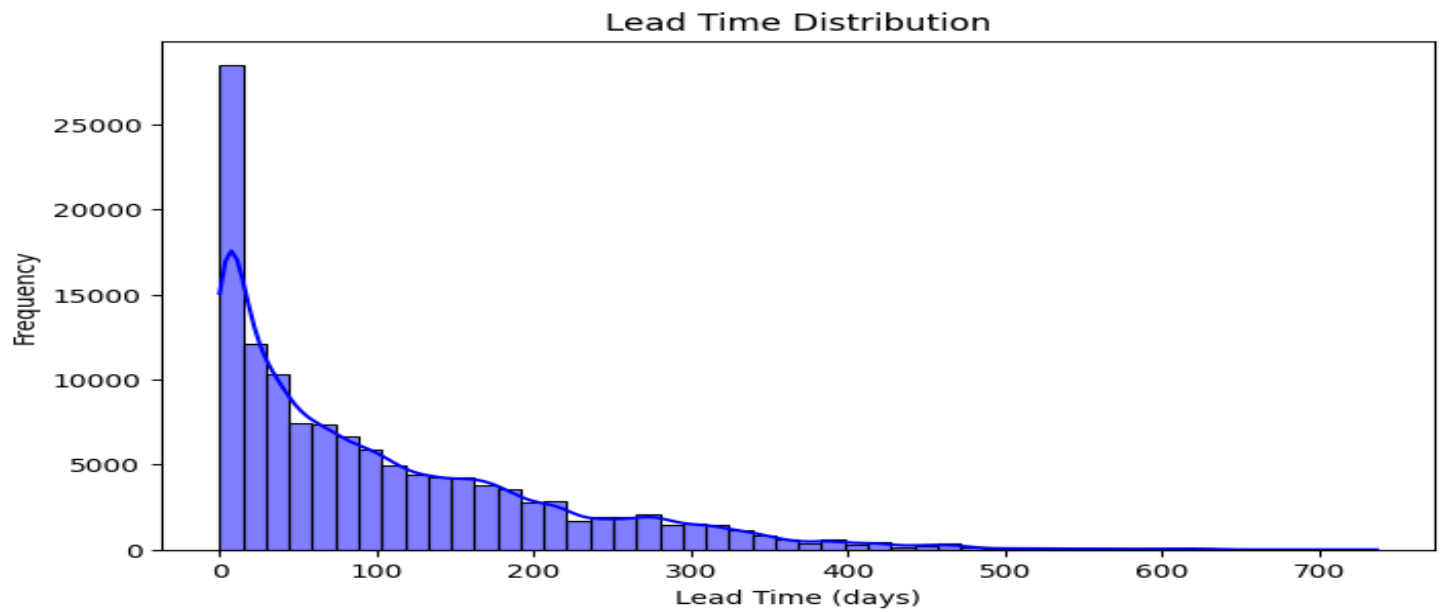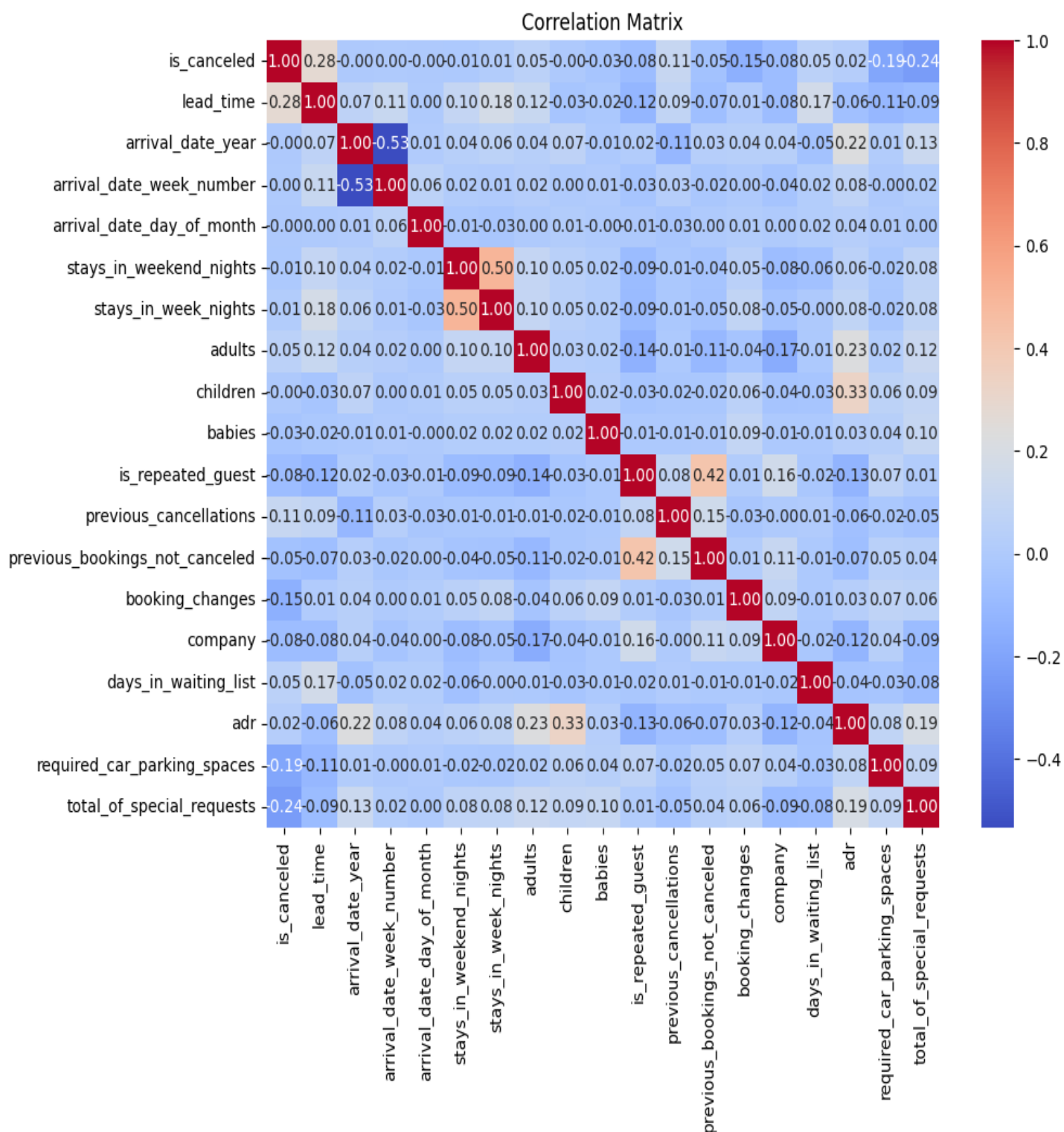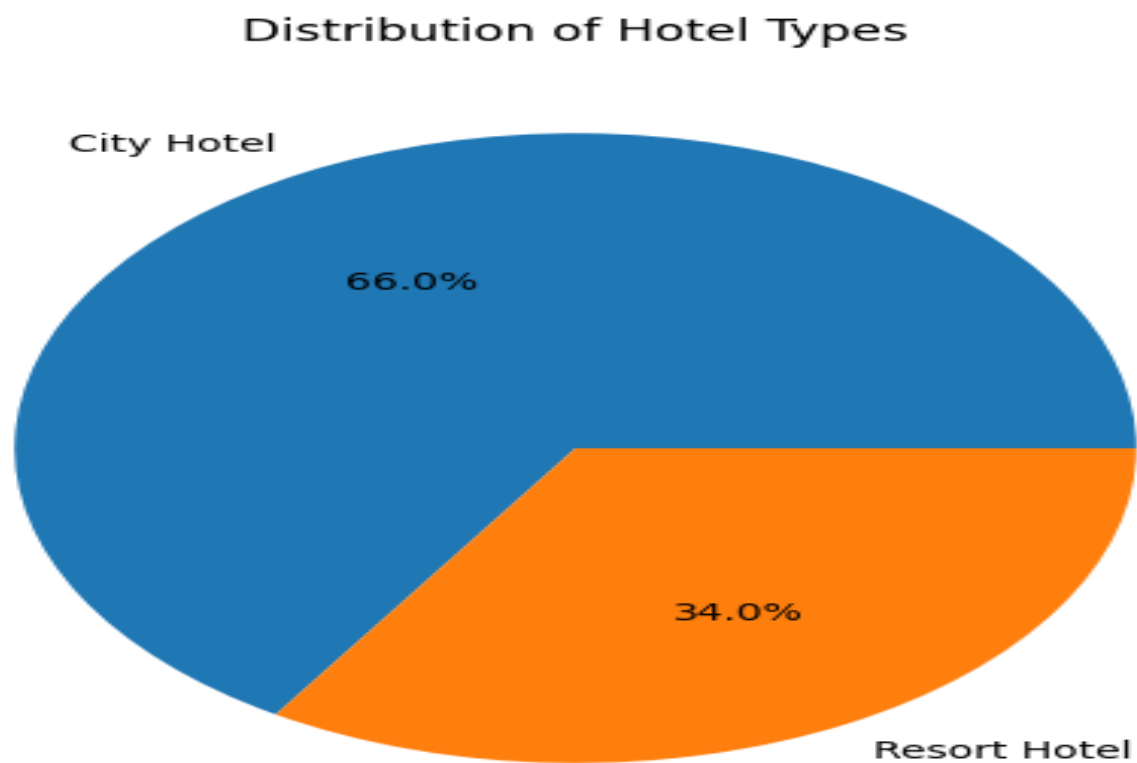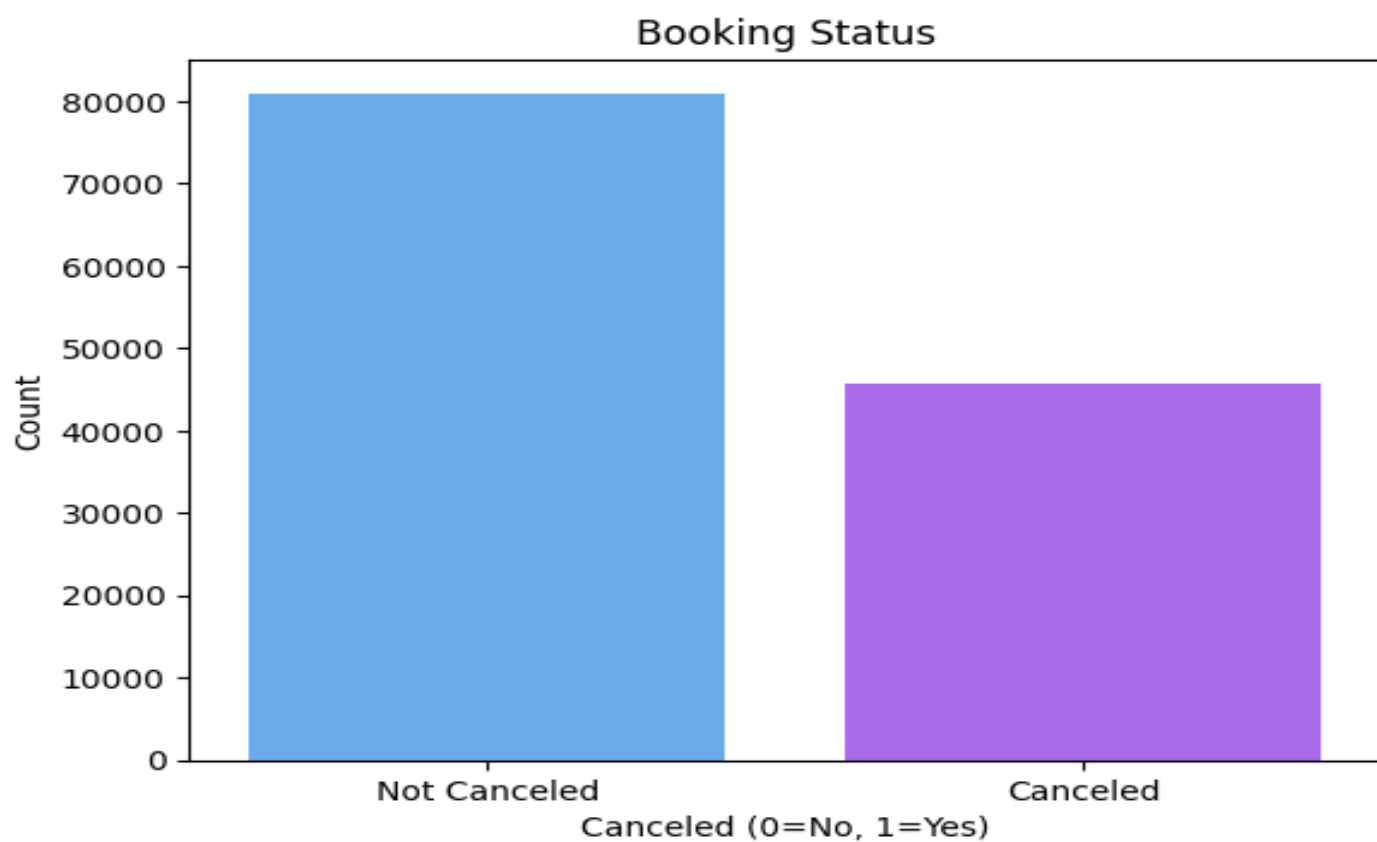
# OUTPUT

## Lead Time Distribution



## Top 10 Countries with Most Bookings

Correlation Matrix

Booking Status



Distribution of Hotel Types

# **CONCLUSION**

In conclusion,The hotel booking prediction analysis conducted in this project reveals several valuable insights into customer behavior, booking patterns, and factors influencing cancellations. By performing data preprocessing and visualization, we gained a clearer understanding of how different features—such as lead time, hotel type, country of origin, and number of guests—impact booking outcomes.Key conclusions include:City hotels receive more bookings than resort hotels.The majority of bookings originate from a small group of countries.Higher lead time is often correlated with a higher likelihood of cancellations.Cancellations and confirmed bookings follow identifiable trends, useful for operational planning.Although no machine learning algorithm was directly implemented in this phase, the groundwork has been laid for classification and regression models that can predict cancellations, booking rates, or revenue forecasts. Algorithms like Logistic Regression, Random Forest, or XGBoost can be integrated later based on the features explored.