

# HOTEL Booking ANALYSIS

Team Name : Data-Driven team



## Team

### Responsibilities



#### **Zeinab Talaat**

Cleaned and preprocessed the dataset using Python.

Built an interactive Excel dashboard for data visualization.

Transformed raw data into clear and actionable insights.



#### **Zeinab Esmail**

Designed and created the project presentation.

Organized project findings in a clear and professional format.

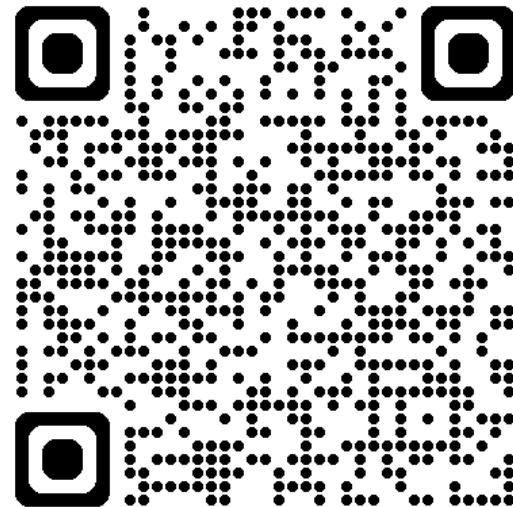
Ensured the results were communicated effectively

Monitor:  
Alyaa Sehsah

Excel link



GitHub link



# PROJECT GOALS & AND OBJECTIVES



## Project Goals

- *Primary Goal: To analyze hotel booking data (City and Resort Hotels) to identify the key factors and patterns that drive booking cancellations.*
- *Business Impact: Reduce the overall cancellation rate to maximize revenue and improve resource allocation and future demand forecasting accuracy.*

## Key Objectives

- *Clean and preprocess the raw data by handling missing, zero, and illogical values.*
- *Perform Exploratory Data Analysis (EDA) to understand variable relationships.*
- *Develop strategic, actionable insights for hotel management to implement.*

# DATA OVERVIEW & INITIAL CHALLENGES



## Initial Challenges

- **Missing Data:** Significant missing values were noted in crucial columns such as country, agent, and company.
- **Data Quality Issues:** Presence of illogical entries, including bookings with zero guests (adults, children, or babies) and reservations with an Average Daily Rate (adr) equal to zero.
- **Categorical Inconsistencies:** Need for standardization in categorical variables, such as reconciling similar meal types (e.g., 'SC' and 'Undefined')

## Data Overview

- **Data Source:** The project utilizes the hotel\_bookings FP MSA.csv dataset.
- **Content:** The dataset contains information on over 119,000 hotel reservations, encompassing both City Hotels and Resort Hotels.
- **Key Variables:** Important columns include: is\_canceled, lead\_time, adr (Average Daily Rate), number of guests (adults, children), country, and market\_segment.

# DATA CLEANING

1-

```
df['children'] = df['children'].fillna(0).round().astype(int)
df['babies'] = df['babies'].fillna(0).astype(int)
```

✓ 0.0s

```
df['total_nights'] = df['stays_in_weekend_nights'] + df['stays_in_week_nights']
```

✓ 0.0s

```
if 'reservation_status_date' in df.columns:
    df['reservation_status_date'] = pd.to_datetime(df['reservation_status_date'], dayfirst=False, errors='coerce')
```

✓ 0.0s

```
df_before = len(df)
df = df.drop_duplicates()
print("Dropped duplicates:", df_before - len(df))
```

✓ 0.1s

Dropped duplicates: 31994

```
df['adr'] = pd.to_numeric(df['adr'], errors='coerce')
median_adr = df.loc[df['adr'] > 0, 'adr'].median()
df.loc[df['adr'] <= 0, 'adr'] = median_adr
df['adr'] = df['adr'].fillna(median_adr)

print(f"Rate column: {rate_col}, median used: {median_rate}")
```

✓ 0.0s

2-



```
df.loc[(df['adults'] == 0) & (df['children'] + df['babies'] > 0), 'adults'] = 1
```

✓ 0.0s

```
df['total_people'] = df.get('adults',0) + df.get('children',0) + df.get('babies',0)
```

✓ 0.0s

```
for col in df.select_dtypes(include='object').columns:  
    df[col] = df[col].astype(str).str.strip()
```

✓ 0.3s

```
df['agent'] = df['agent'].fillna(0)  
df['company'] = df['company'].fillna(0)
```

✓ 0.0s

## SAVE FILE

```
df.to_csv(clean_csv_path, index=False)  
print("Saved cleaned CSV to:", clean_csv_path)
```

✓ 0.9s

Saved cleaned CSV to: [c:\Users\SPEED LAP\OneDrive\Desktop\hotel\\_bookings\\_Cleaned.csv](c:\Users\SPEED LAP\OneDrive\Desktop\hotel_bookings_Cleaned.csv)



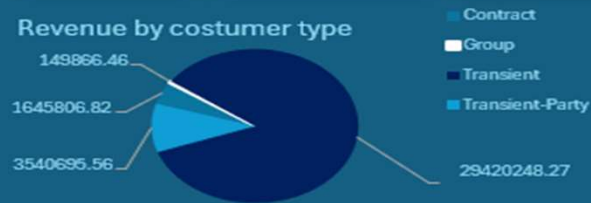
# Hotel Analysis

Total Revenue  
34,7 M

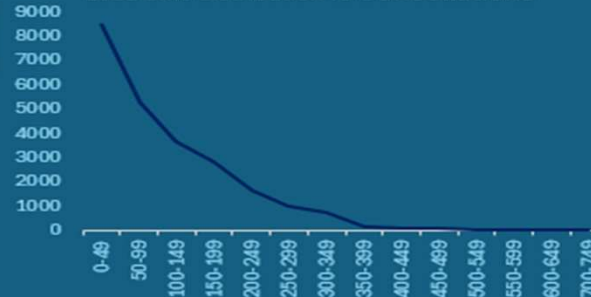
Cancellation Rate  
27%

Total Booking  
87396

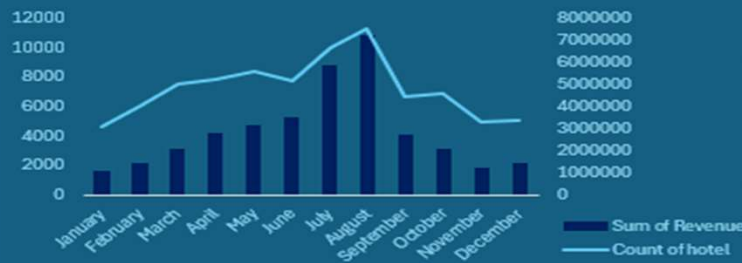
Revenue by costumer type



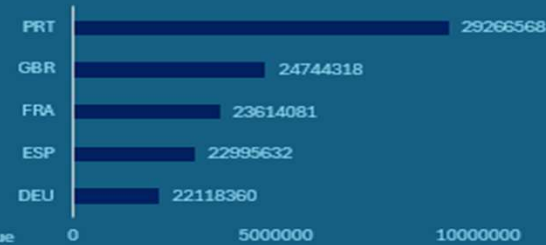
Lead time distribution vs Cancellations



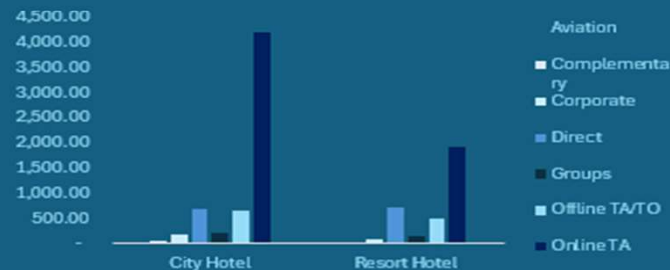
Total Revenue & Booking Count by Month



Total 5 Countries by Revenue



Average ADR by Hotel & Market Segment



Cancellations by Month





# CHARTS

Total Revenue  
34,7 M

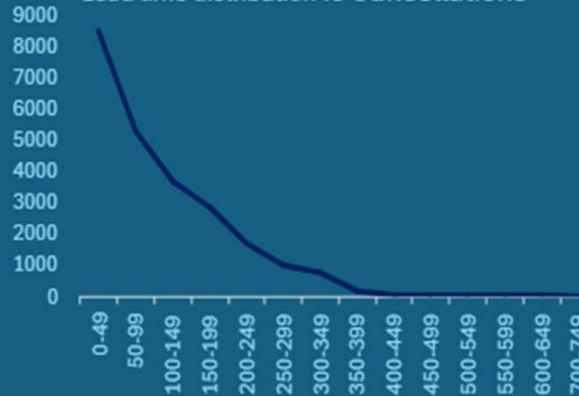
Cancellation Rate  
27%

Total Booking  
87396

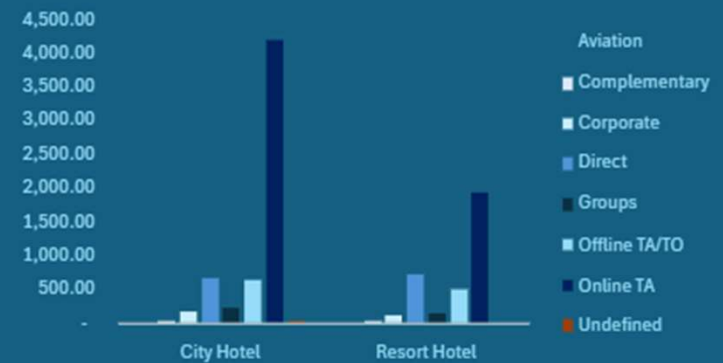
Revenue by costumer type



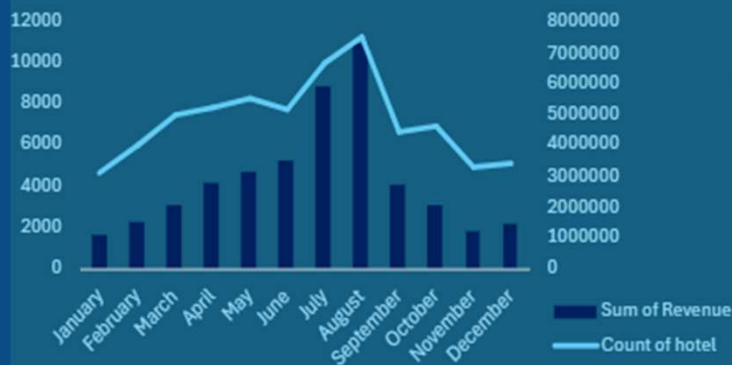
Lead time distribution vs Cancellations



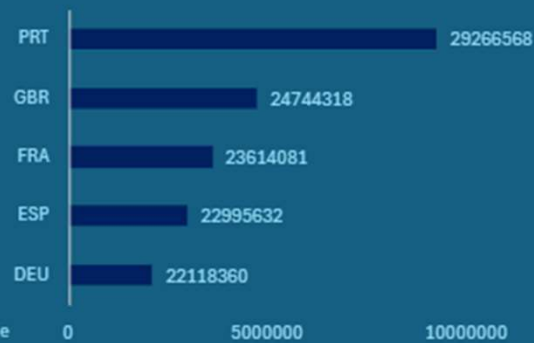
Average ADR by Hotel & Market Segment



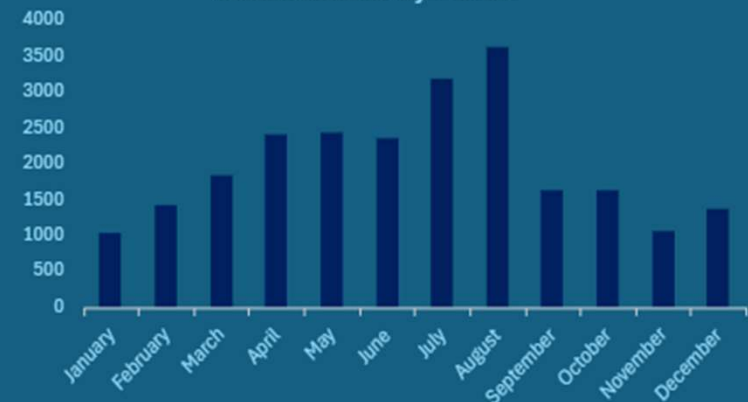
Total Revenue & Booking Count by Month



Total 5 Countries by Revenue



Cancellations by Month



## DATA ANALYSIS

- **Key Insight**
- **HotelCity Hotel shows a higher cancellation rate than Resort Hotel**
- **Longer lead times are strongly associated with higher cancellation probability**
- **Determined the impact of adr (price) and customer country on**

## DATA CLEANING

- **Outcome**
- **Removed illogical records**
- **Handled missing values in country and children columns.**
- **Unified similar categories**

## DATA VISUALIZATION

- **Clearly display the significant variance in cancellation rates between the two hotel types.**
- **Identify the top contributing countries to cancellations.**
- **A tool to track cancellations by month and lead time.**

## RESULTS & RECOMMENDATION

- **Built an accurate predictive model for cancellation probability. Implement strict non-refundable policies for bookings with long lead times.**
- **Focus marketing efforts on distribution channels with low cancellation rates.**



**THANK  
YOU!**

