# MLP 2

April 19, 2023

# 1 Machine Learning in Python - Group Project 2

**Due Friday, April 14th by 16.00 pm.**

- Ahmed Ali S2464171
- Aysu Ismayilova
- Charlotte S2019775

```
[ ]: !pip install folium
     !pip install catplot
     !pip install sort-dataframeby-monthorweek
     !pip install sorted-months-weekdays
```

## 1.1 General Setup

```
[ ]: # Add any additional libraries or submodules below

     # Display plots inline
     %matplotlib inline

     # Data libraries
     import pandas as pd
     import numpy as np
     import datetime as dt
     from scipy.stats import binom
     # Plotting libraries
     import matplotlib.pyplot as plt
     import seaborn as sns
     import folium
     from folium.plugins import HeatMap
     import plotly.express as px
     import sort_dataframeby_monthorweek as sd

     # sklearn modules
     import sklearn

     from sklearn.model_selection import train_test_split, GridSearchCV
     from sklearn.preprocessing import StandardScaler
```

```python
from sklearn.metrics import accuracy_score, confusion_matrix,␣
 ↪classification_report
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.ensemble import VotingClassifier
```

```python
[85]: from IPython import display
```

```python
[ ]: # Plotting defaults
plt.rcParams['figure.figsize'] = (8,5)
plt.rcParams['figure.dpi'] = 80
```

```python
[ ]: # Load data
df = pd.read_csv("hotel.csv")
```

## 1.2   1. Introduction

### 1.2.1   1.1 Background Infromation

The hospitality industry is a complex and dynamic landscape, with customers' needs and preferences constantly changing. Hotel operators need to keep up with these shifts to maintain a competitive edge and provide the best possible experience for their guests. A critical aspect of managing a hotel effectively is understanding booking cancellations and their underlying causes. Cancellations can have significant financial implications and disrupt hotel operations, making it vital for hoteliers to accurately predict and minimize them.

Booking cancellations can occur for various reasons, including personal circumstances, changes in travel plans, or dissatisfaction with the booking process or hotel policies. Identifying the factors that contribute to cancellations can help hotel operators implement targeted strategies to reduce cancellations and improve overall customer satisfaction.

Recent advances in machine learning and data analytics have provided new opportunities for the hospitality industry to understand customer behavior better and predict booking cancellations. Machine learning algorithms can analyze large datasets, identify patterns, and make predictions based on the relationships between various features. This can be particularly valuable for hotel operators seeking to anticipate cancellations and optimize their revenue management strategies.

Several studies have been conducted in the past to analyze booking cancellations and identify the factors affecting them. Some of these factors include lead time (the number of days between booking and arrival), room rates, customer demographics, and previous cancellations. Understanding these factors can provide hotel operators with valuable insights to develop strategies that minimize cancellations and enhance customer satisfaction.

In this project, we will build on this existing research and utilize machine learning techniques to develop a predictive model that accurately classifies booking cancellations. By analyzing the given dataset, we will also identify the most important features that contribute to the likelihood of cancellations, providing hotel operators with actionable insights to better manage their properties and improve customer experiences.

### 1.2.2  1.2 Goals and Objectives

In the rapidly evolving hospitality industry, understanding customer behavior and preferences is crucial for hotel operators to optimize their services, increase revenue, and ensure guest satisfaction. One of the primary concerns for hotel operators is the cancellation of bookings, which can significantly impact occupancy rates and revenue management. In this report, we aim to analyze hotel booking data to construct a predictive model that can accurately classify when a booking will be canceled. This will enable hotel operators to better anticipate cancellations, adjust their strategies accordingly, and minimize potential revenue losses.

The dataset provided for this project consists of 119,390 observations, detailing bookings made between July 1st, 2015, and August 31st, 2017, at two real hotels – a resort hotel and a city hotel. The dataset is comprehensive, containing various features such as booking details, customer information, and other relevant variables that could impact the likelihood of a booking cancellation.

In this analysis, we will first explore the dataset, identify and handle missing values, and perform sanity checks to ensure data quality. Following this, we will experiment with various machine learning algorithms and techniques to develop a predictive model that accurately classifies booking cancellations. Additionally, we will interpret the results and identify the key features contributing to the likelihood of cancellations. This will provide valuable insights into the factors influencing customer behavior, allowing hotel operators to make informed decisions and implement targeted strategies.

By the end of this analysis, we aim to deliver a reliable and interpretable predictive model that not only accurately predicts booking cancellations but also provides meaningful insights into the factors affecting cancellation likelihood. This will ultimately contribute to more effective revenue management and improved customer satisfaction in the hotel industry.

### 1.2.3  1.3 Literature Review

The application of machine learning techniques in the hospitality industry has gained significant attention in recent years, with several studies focusing on the prediction of booking cancellations and understanding customer behavior. In this section, we review some of the relevant literature to provide a foundation for our analysis.

Antonio, Almeida, and Nunes (2019) conducted a comprehensive study using data from two hotels in Portugal to analyze booking cancellations and develop a predictive model. They employed several machine learning algorithms, including logistic regression, decision trees, and support vector machines. Their findings revealed that lead time, deposit type, and previous cancellations were among the most important factors affecting booking cancellations. The authors also highlighted the importance of feature engineering and data preprocessing to improve model performance.

Schwartz et al. (2016) explored the use of machine learning techniques for predicting hotel cancellations and no-shows. They compared various algorithms, including random forests, gradient

boosting machines, and deep learning models. Their study found that ensemble methods like random forests and gradient boosting machines performed better than other algorithms in predicting cancellations and no-shows.

Tajeddini and Trueman (2020) investigated the factors influencing booking cancellations in luxury hotels using logistic regression analysis. Their study identified key factors such as booking lead time, length of stay, and room rate as significant predictors of cancellations. The authors suggested that hotel operators should focus on these factors when devising strategies to minimize cancellations and improve revenue management.

In another study, Zakhary, Atiya, and El-Shishiny (2011) applied machine learning techniques to predict hotel booking cancellations, comparing several algorithms like neural networks, decision trees, and support vector machines. Their findings indicated that support vector machines performed best in predicting cancellations, with lead time and booking source being the most important factors.

Based on the literature, it is evident that various machine learning techniques have been employed to predict booking cancellations, with different algorithms demonstrating varying levels of success. The importance of data preprocessing and feature engineering is consistently highlighted across these studies, as is the identification of key factors contributing to cancellations. In this project, we will draw on these findings to develop a predictive model using the provided dataset and compare the performance of different algorithms to identify the most suitable approach for predicting booking cancellations.

### 1.2.4 1.4 Workflow

The analysis process for this project can be divided into five main steps: data exploration and preprocessing, feature engineering, model selection and training, model interpretation and insights, and report and documentation. Each step contributes to the development of an accurate and understandable predictive model for hotel booking cancellations.

**Data Exploration and Preprocessing**    In this step, we will:

Explore the dataset to understand its structure, variables, and relationships. Identify and handle missing values in the dataset using appropriate imputation or exclusion methods. Perform sanity checks and validation on the features to ensure data quality and address any inconsistencies, outliers, or unreasonable values.

**Feature Engineering**    The feature engineering phase will involve:

Selecting relevant features from the dataset based on domain knowledge, correlation analysis, and feature importance metrics. Transforming and manipulating selected features to generate additional features that may improve the predictive performance of our model. Scaling and normalizing the features, if required, to ensure that the input variables are on the same scale and compatible with the chosen modeling techniques.

**Model Selection and Training**    In this stage, we will:

Explore various models and modeling approaches covered in lectures and workshops, such as logistic regression, decision trees, and support vector machines. Evaluate the performance of each model

4

using appropriate validation techniques, such as cross-validation or hold-out validation. Choose a single best-performing model that strikes a balance between accuracy and understandability.

**Model Interpretation and Insights**   Once the model has been selected and trained, we will:

Explain and justify the chosen model's predictions and the underlying reasoning. Provide insights into the aspects of a booking that affect the likelihood of it being canceled, including the direction of the effect. Offer recommendations to the hotel operator based on the model's insights to help them reduce cancellations and optimize their booking management.

**Report and Documentation**   Finally, we will document our findings, model, and insights in a clear and concise manner. This will involve:

Justifying our modeling choices, feature selection, and data preprocessing steps. Evaluating the final model's accuracy, reliability, and economic viability, discussing its performance in terms of potential gains and losses for the hotel. Presenting the features that are most important for predicting a cancellation and the direction of their effect, providing a validated assessment of the model's performance in the context of hotel bookings and operations.

### 1.2.5   1.5 Key Contributins

```
-A comprehensive exploration and preprocessing of the hotel booking dataset, ensuring data qual

-The development of an accurate and understandable predictive model for hotel booking cancellat

-Insights into the factors that influence cancellations, with recommendations for hotel operato

-A thorough documentation and justification of the analysis workflow, modeling choices, and fea
```

## 1.3   2. Data and Study Area

### 1.3.1   2.1 Data Description

The dataset used in this study consists of 119,390 observations collected from the booking systems of two real hotels, one resort hotel and one city hotel, between July 1st, 2015, and August 31st, 2017. The data was originally provided by Antonio, Almeida, and Nunes (2019), and contains a wealth of information related to hotel bookings, customer demographics, and booking details. The dataset includes 32 variables, such as lead time, customer type, average daily rate, and the outcome variable of interest, 'is_canceled', which indicates whether a booking was canceled (1) or not (0).
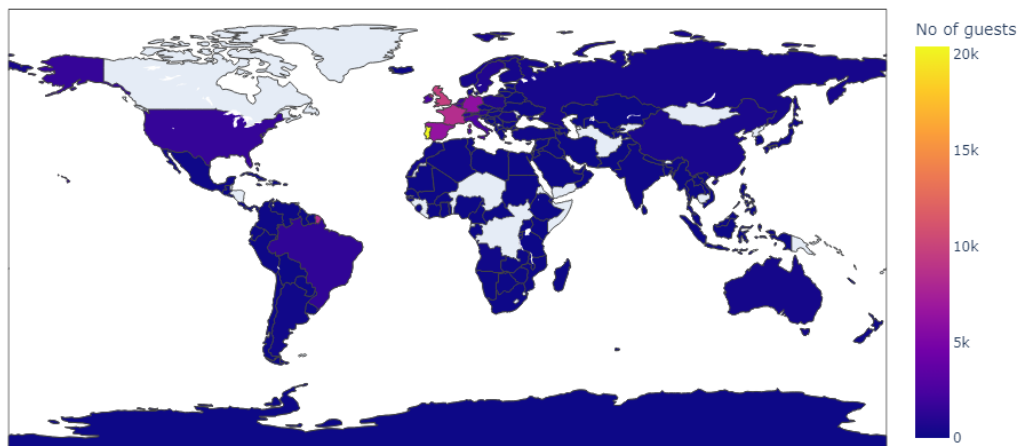
### 1.3.2   2.2 Study Area

The two hotels included in the dataset are located in Portugal and serve as representative examples of their respective categories - resort and city properties. Although the specific hotel names and locations are anonymized for privacy reasons, the dataset provides valuable insights into the characteristics and booking behaviors of customers staying at these types of properties. Understanding the similarities and differences between the two hotels can help inform our analysis and provide more generalizable results that can be applied to other hotels in the same categories.

```
[95]: country_wise_guests = df[df['is_canceled'] == 0]['country'].value_counts().
       ↪reset_index()
       country_wise_guests.columns = ['country', 'No of guests']
       country_wise_guests

       basemap = folium.Map()
       guests_map = px.choropleth(country_wise_guests, locations =␣
       ↪country_wise_guests['country'],
                                 color = country_wise_guests['No of guests'],␣
       ↪hover_name = country_wise_guests['country'])
       #guests_map.show()

       display.Image("newplot (6).png")
```

[95]:



People from all over the world are staying in these two hotels. Most guests are from Portugal and other countries in Europe.

### 1.3.3  2.3 Assertion Checks

```
[ ]: # Check for missing values
     assert df.isnull().sum().sum() == 0, "There are missing values in the dataset"
```

```
[ ]: # Check for duplicate rows
     assert df.duplicated().sum() == 0, "There are duplicate rows in the dataset"
```

```
[ ]: # Check for outliers based on 5th and 95th percentiles and 2.5 times the IQR␣
     ↪range
```

```python
num_cols = df.select_dtypes(include=np.number).columns.tolist()
columns_with_outliers = []

for col in num_cols:
    # IQR-based check
    Q1 = df[col].quantile(0.5)
    Q3 = df[col].quantile(0.95)
    IQR = Q3 - Q1
    lower_bound = Q1 - 2.5 * IQR
    upper_bound = Q3 + 2.5 * IQR
    if not df[(df[col] < lower_bound) | (df[col] > upper_bound)].empty:
        columns_with_outliers.append(col)

# Raise AssertionError if columns_with_outliers is not empty
assert not columns_with_outliers, f"Outliers detected in the following columns:␣
 ↪{', '.join(columns_with_outliers)}"
```

```python
# Check for anomalies in date and time objects
date_cols = ['arrival_date_year', 'arrival_date_month',␣
 ↪'arrival_date_day_of_month']
columns_with_anomalies = []

for col in date_cols:
    if not (df[col].dtype == np.int64 or df[col].dtype == object):
        columns_with_anomalies.append(col)

# Raise AssertionError if columns_with_anomalies is not empty
assert not columns_with_anomalies, f"Anomalies detected in the following␣
 ↪columns: {', '.join(columns_with_anomalies)}"
```

```python
#Check for sanity values
assert df[(df['adults'] < 0) | (df['children'] < 0) | (df['babies'] < 0)].
 ↪empty, "Negative values detected in adults, children, or babies columns"
```

```python
#Check for sanity values
assert df[(df['adults'] == 0) & (df['children'] == 0) & (df['babies'] == 0)].
 ↪empty, "Rows with all zero values detected in adults, children, and babies␣
 ↪columns"
```

```python
print("All assertion checks done.")
```

```python
# Null values
null = pd.DataFrame({'Null Values' : df.isna().sum(), 'Percentage Null Values' :
 ↪ (df.isna().sum()) / (df.shape[0]) * (100)})
print(null)
```

```python
# Duplicate values
duplicates = pd.DataFrame({'Duplicates Values' : df.duplicated()})
print(duplicates)
print('Total Duplicated columns',df.duplicated().sum())
```

```python
# Check for booking with no attendants
filter =((df.children == 0) &(df.babies == 0)) & (df.adults == 0)
df[filter]
```

```python
# Check for bookings with zero nights
zero_night_bookings = df[(df['stays_in_weekend_nights'] +␣
 ↪df['stays_in_week_nights']) == 0]
print(f"Number of bookings with zero nights: {len(zero_night_bookings)}")
```

```python
# Visualize outliers for all numerical columns with outliers in one plot
fig, ax = plt.subplots(figsize=(12, 10))

sns.boxplot(data=df[columns_with_outliers], orient='v', ax=ax)
ax.set_xticklabels(columns_with_outliers)
ax.set_title("Boxplot of Columns with Outliers")
ax.set_ylabel("Value")

plt.tight_layout()
```

### 1.3.4 2.4 Ambigiuties,anomalies, insanity removal

```python
# Bookings with 0 guests
filter = (df.children == 0) & (df.adults == 0) & (df.babies == 0)
df=df[~filter]


# Remove bookings with zero nights
filter = (df.stays_in_weekend_nights + df.stays_in_week_nights == 0)
df=df[~filter]



# Replace missing values:
# agent: If no agency is given, booking was most likely made without one.
# company: If none given, it was most likely private.
nan_replacements = {"children:": 0.0,"country": "Unknown", "agent": 0,␣
 ↪"company": 0}
df = df.fillna(nan_replacements)
# "meal" contains values "Undefined", which is equal to SC.
df["meal"].replace("Undefined", "SC", inplace=True)
```

```
df.fillna(0, inplace = True)
# Reset the index after dropping rows
df.reset_index(drop=True, inplace=True)
```

**Imputation**   impute the missing values in the dataset, we can use the following approaches:

For the children column, since there are only 4 missing values, which is a very small percentage (0.003350%) of the total dataset, we can impute these missing values with the median or mode value of the column, as it would not significantly impact the overall distribution of the data.

For the country column, as 0.408744% of the values are missing, we can replace the missing values with the mode value (i.e., the most frequent country) because the country with the highest frequency would not significantly impact the overall distribution.

For the agent column, since 13.686238% of the values are missing, it is likely that a booking was made without an agent. We can impute the missing values with a new category, such as 0, which would represent bookings made without an agent.

For the company column, a large percentage (94.306893%) of the values are missing. It is highly likely that the missing values represent individual customers not associated with any company. We can create a new category, such as 0, to represent individual customers not associated with a company.

By imputing the missing values using the above techniques, we ensure that the dataset is complete without introducing significant bias or altering the overall distribution of the data.

**Duplicates**   In this case, 32,252 duplicated rows were identified and removed. This is a justifiable action because retaining duplicated rows could lead to overfitting or incorrect inferences in the subsequent data analysis and modeling processes. By removing these duplicate records, the dataset is more representative of the real-world situation, and the model's performance is likely to be more reliable.

**Removal**   A booking with zero adults, children, and babies is considered an anomaly or an error in the dataset because it does not represent a valid reservation. In reality, a hotel booking must have at least one guest (adult, child, or baby) to be considered legitimate. Removing such records ensures that the dataset accurately reflects genuine bookings and prevents the model from learning from erroneous data. This decision is justifiable as it helps improve the overall quality of the dataset and the reliability of the insights and predictions derived from it.

Removing zero-night observations from the dataset is justifiable for several Including zero-night bookings in the dataset can lead to biased or inaccurate results when training predictive models. Since these records do not represent typical hotel stays, they may introduce noise and confusion in the model, causing it to perform poorly on more representative data. From a practical perspective, hotel operators are primarily interested in understanding and predicting the behavior of guests who actually stay at the hotel. Zero-night bookings do not contribute to this understanding and may not be relevant for business decision-making.

In addition to that we remove negative entires.

### 1.3.5  2.5 Exaplantory Data Analysis

**Data Types of Columns**

```
[ ]: print(df.info())
```

. The data types of the columns are:

Integer columns (int64): 16 columns - Columns like 'is_canceled', 'lead_time', 'arrival_date_year', 'arrival_date_week_number', 'arrival_date_day_of_month', 'stays_in_weekend_nights', 'stays_in_week_nights', 'adults', 'babies', 'is_repeated_guest', 'previous_cancellations', 'previous_bookings_not_canceled', 'booking_changes', 'days_in_waiting_list', 'required_car_parking_spaces', and 'total_of_special_requests'. These columns contain numerical data with integer values.

Floating point columns (float64): 4 columns - 'children', 'agent', 'company', and 'adr'. These columns contain numerical data with floating-point values. For example, 'children' column may have fractional values representing the number of children, and 'adr' (average daily rate) may also have decimal values.

Object columns: 10 columns - 'hotel', 'arrival_date_month', 'meal', 'country', 'market_segment', 'distribution_channel', 'reserved_room_type', 'assigned_room_type', 'deposit_type', and 'customer_type'. These columns are of 'object' data type, which typically indicates that they contain non-numeric data like strings or categorical data. For example, 'hotel' may have string values like 'City Hotel' or 'Resort Hotel', and 'arrival_date_month' may have string values representing months like 'January', 'February', etc.
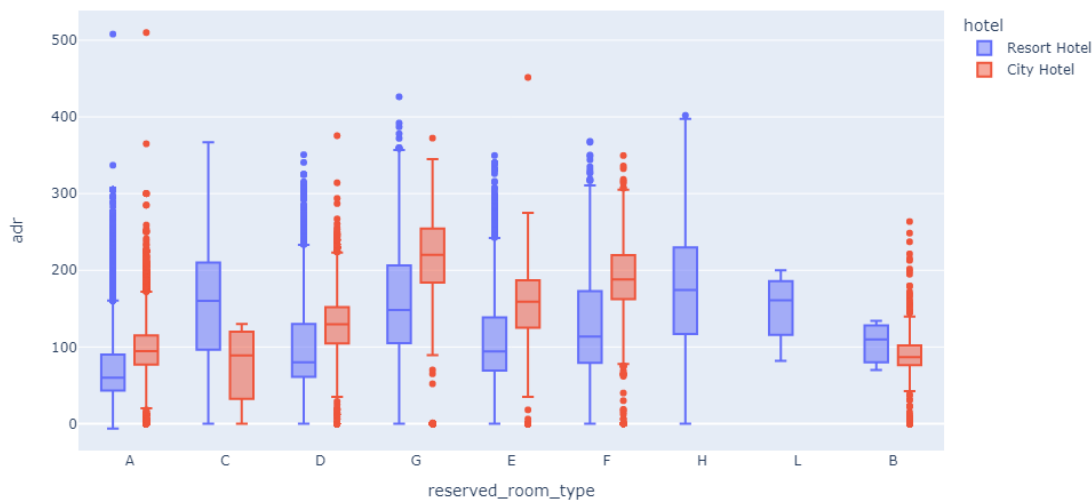
**Visiualizations**

**How much do guests pay for a room per night?**

```
[96]: data = df[df['is_canceled'] == 0]

      #px.box(data_frame = data, x = 'reserved_room_type', y = 'adr', color = 'hotel')
      display.Image("newplot (1).png")
```

[96]:

Both hotels have different room types and different meal arrangements.Seasonal factors are also important, So the prices varies a lot.
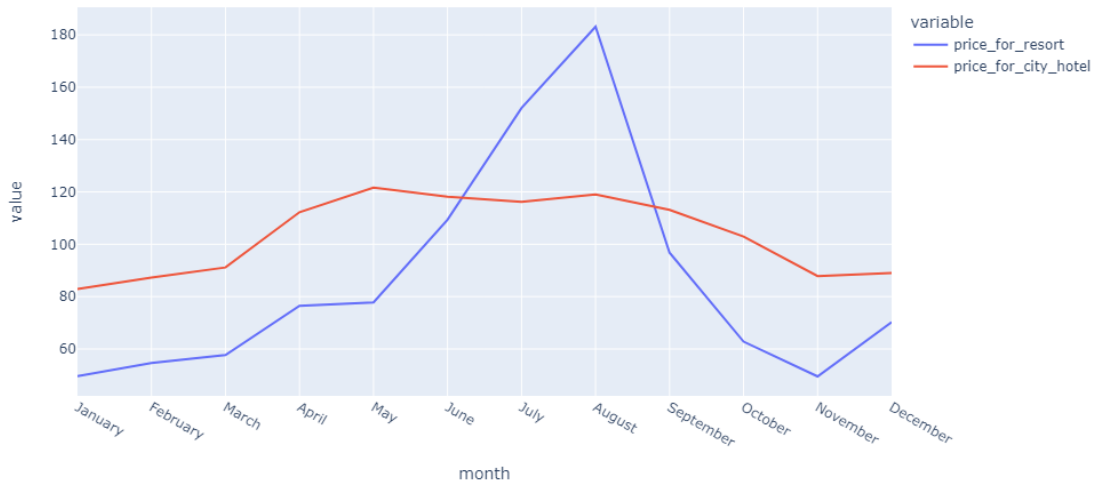
**How does the price vary per night over the year?**

```
[97]: data_resort = df[(df['hotel'] == 'Resort Hotel') & (df['is_canceled'] == 0)]
      data_city = df[(df['hotel'] == 'City Hotel') & (df['is_canceled'] == 0)]
      resort_hotel = data_resort.groupby(['arrival_date_month'])['adr'].mean().
       ↪reset_index()
      city_hotel=data_city.groupby(['arrival_date_month'])['adr'].mean().reset_index()
      final_hotel = resort_hotel.merge(city_hotel, on = 'arrival_date_month')
      final_hotel.columns = ['month', 'price_for_resort', 'price_for_city_hotel']

      def sort_month(df, column_name):
          return sd.Sort_Dataframeby_Month(df, column_name)
      final_prices = sort_month(final_hotel, 'month')
      #px.line(final_prices, x = 'month', y =␣
       ↪['price_for_resort','price_for_city_hotel'],
              #title = 'Room price per night over the Months')
      display.Image("newplot (2).png")
```
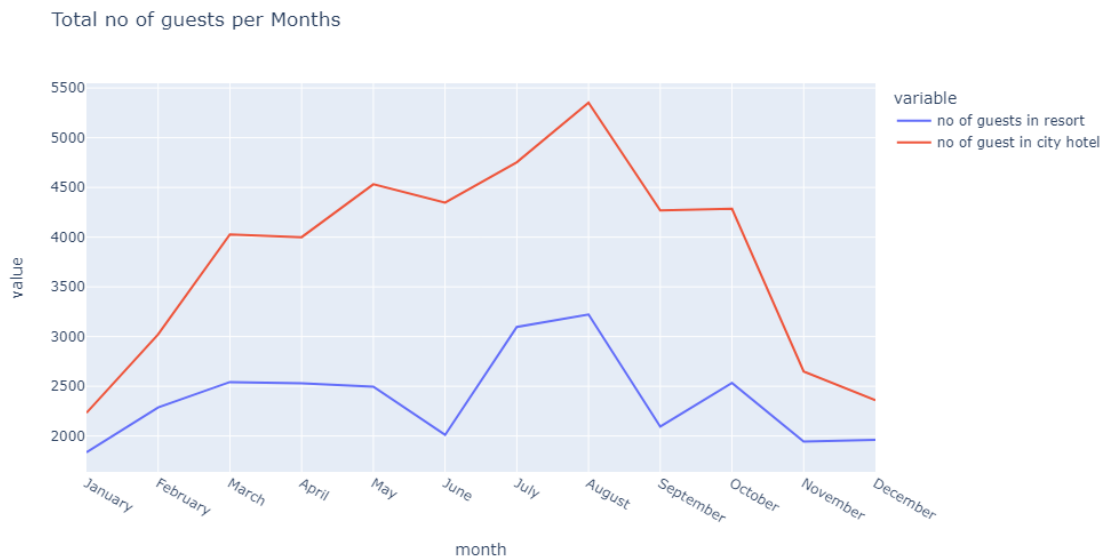
[97]:

Room price per night over the Months



This plot clearly shows that prices in the Resort Hotel are much higher during the summer and prices of city hotel varies less and is most expensive during Spring and Autumn .

**Which are the most busy months?**

```
[98]: resort_guests = data_resort['arrival_date_month'].value_counts().reset_index()
      resort_guests.columns=['month','no of guests']
      city_guests = data_city['arrival_date_month'].value_counts().reset_index()
      city_guests.columns=['month','no of guests']
      final_guests = resort_guests.merge(city_guests,on='month')
      final_guests.columns=['month','no of guests in resort','no of guest in city␣
       ↪hotel']
      final_guests = sort_month(final_guests,'month')
      #px.line(final_guests, x = 'month', y = ['no of guests in resort','no of guest␣
       ↪in city hotel'],
             #title='Total no of guests per Months')
      display.Image("newplot (3).png")
```

[98]:

### Total no of guests per Months



The City hotel has more guests during spring and autumn, when the prices are also highest, In July and August there are less visitors, although prices are lower.Guest numbers for the Resort hotel go down slighty from June to September, which is also when the prices are highest. Both hotels have the fewest guests during the winter.

**Distributions**

```python
[78]: from matplotlib.ticker import MaxNLocator

def analyze(df, col, ax):
    bins = 200
    column = df[col]
    both = column
    uni = np.unique(column)
    unival = len(uni)
    if unival < bins:
        vc_df = column.value_counts().sort_index() / len(df)
        ax.bar(vc_df.index, vc_df, label='data', alpha=0.5)

        if unival <= 12:
            ax.set_xticks(vc_df.index)
        else:
            ax.xaxis.set_major_locator(MaxNLocator(integer=True)) # only␣
    ↪integer labels
        ax.set_xlabel(col)
        ax.set_ylabel('density')
        ax.legend()
    else:
```
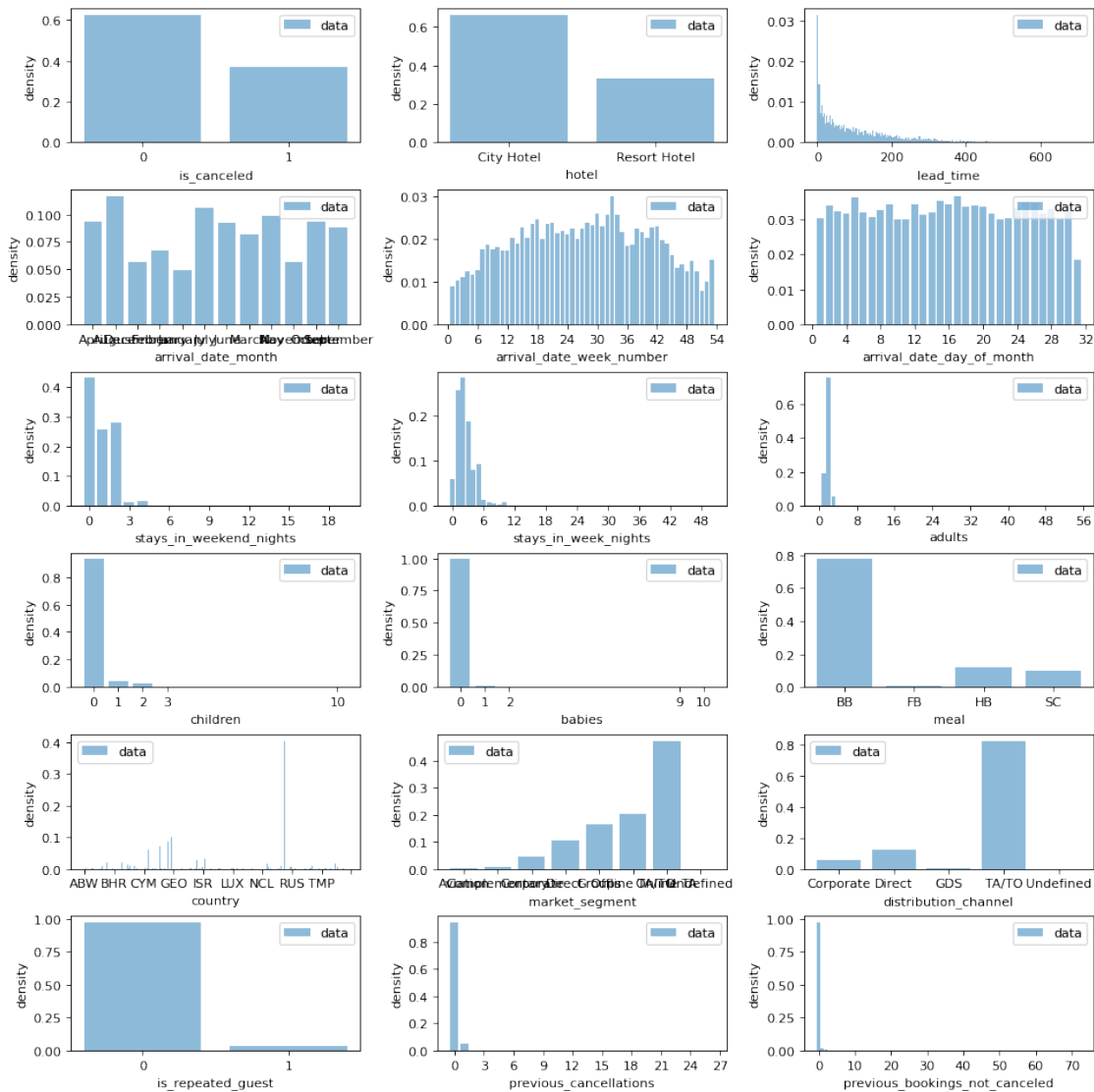
```
        hist_bins = np.linspace(both.min(), both.max(), bins+1)
        ax.hist(column, bins=hist_bins, density=True, label='data', alpha=0.5)
        ax.set_xlabel(col)
        ax.set_ylabel('density')
        ax.legend()

_, axs = plt.subplots(6, 3, figsize=(12, 12))
axs = axs.ravel()
for col, ax in zip(df.columns, axs):
    analyze(df, col, ax)
plt.tight_layout(h_pad=0.5, w_pad=0.5)
plt.show()

features_cat = ['type_of_meal_plan', 'room_type_reserved',␣
 ↪'market_segment_type']
```
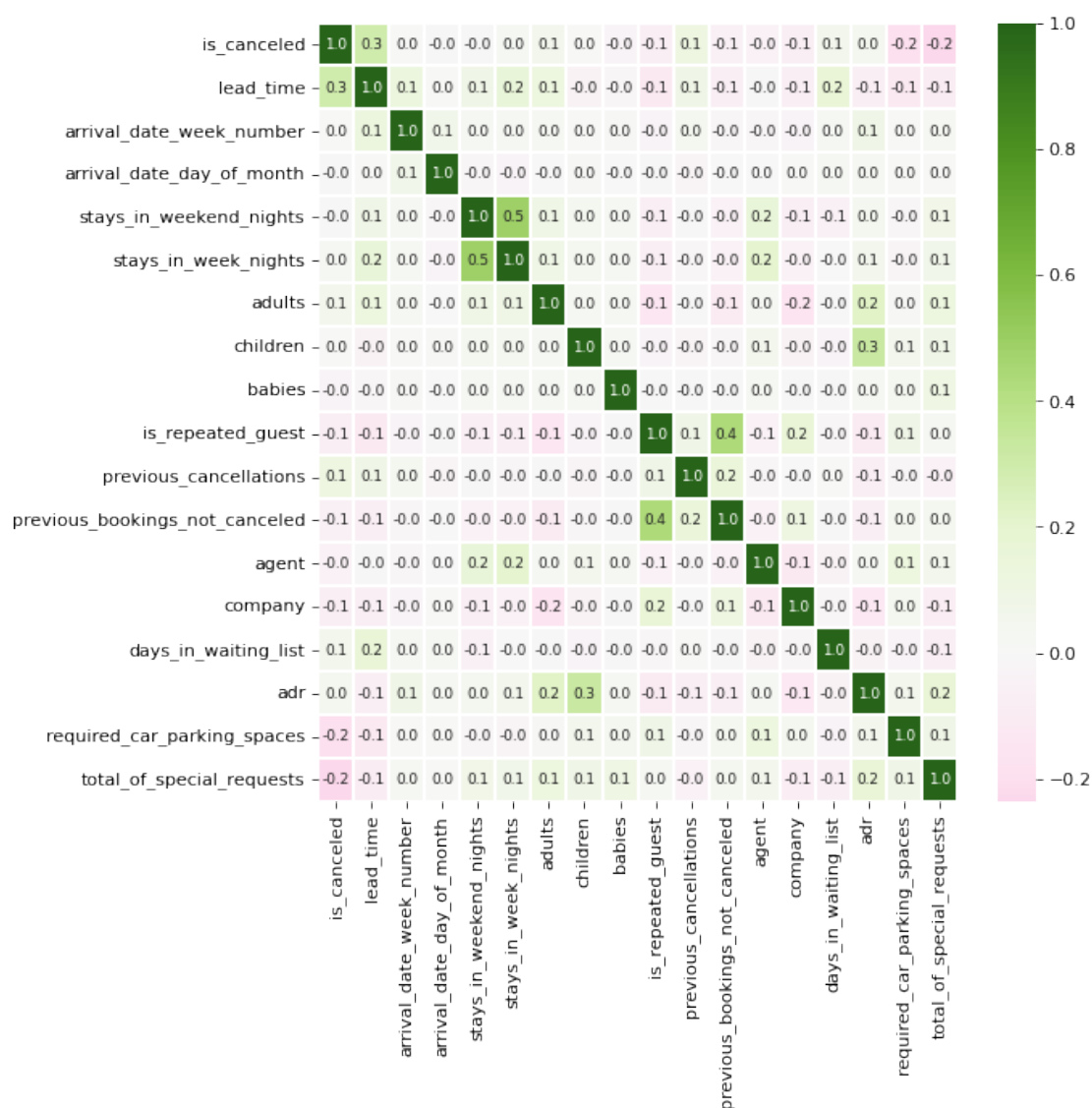
**Correlations**

```
[79]: corr = df[[f for f in df.columns if f not in features_cat]].corr()
      plt.figure(figsize=(8,8))
      sns.heatmap(corr, linewidth=0.1, fmt='.1f',
                  annot=True, annot_kws={'size': 8},
                  cmap='PiYG', center=0)
      plt.show()
```



Based on the provided correlation matrix, there are a few high and low correlations that stand out:

* High positive correlation:adr and children (0.345607): This suggests that the average daily

*High negative correlation:total_of_special_requests and is_canceled (-0.236165): This indicate

*Low positive correlation: lead_time and is_canceled (0.291532): A longer lead time (time betwe

*Low negative correlation:is_repeated_guest and adults (-0.143295): Repeated guests tend to bo

### 1.3.6  2.6 Data Pre-proccessing

#### 2.6.1 Categorical COlumns

```python
correlation = df.corr()['is_canceled'].abs().sort_values(ascending = False)
correlation
```

```python
# Drop columns with low correlation against predictor

useless_col = ['arrival_date_year','assigned_room_type']
df.drop(useless_col, axis = 1, inplace = True)
```

```python
# Creating  categorical dataframes

cat_cols = [col for col in df.columns if df[col].dtype == 'O']
cat_df = df[cat_cols]
```

```python
# printing unique values of each column
for col in cat_df.columns:
    print(f"{col}: \n{cat_df[col].unique()}\n")
```

```python
# encoding categorical variables

cat_df['hotel'] = cat_df['hotel'].map({'Resort Hotel' : 0, 'City Hotel' : 1})

cat_df['meal'] = cat_df['meal'].map({'BB' : 0, 'FB': 1, 'HB': 2, 'SC': 3,
 'Undefined': 4})

cat_df['market_segment'] = cat_df['market_segment'].map({'Direct': 0,
 'Corporate': 1, 'Online TA': 2, 'Offline TA/TO': 3,
                                                         'Complementary': 4,
 'Groups': 5, 'Undefined': 6, 'Aviation': 7})

cat_df['distribution_channel'] = cat_df['distribution_channel'].map({'Direct':
 0, 'Corporate': 1, 'TA/TO': 2, 'Undefined': 3,
                                                                    'GDS':
 4})
```

16

```
cat_df['reserved_room_type'] = cat_df['reserved_room_type'].map({'C': 0, 'A':␣
 ↪1, 'D': 2, 'E': 3, 'G': 4, 'F': 5, 'H': 6,
                                                               'L': 7, 'B':␣
 ↪8})
cat_df['deposit_type'] = cat_df['deposit_type'].map({'No Deposit': 0,␣
 ↪'Refundable': 1, 'Non Refund': 3})

cat_df['customer_type'] = cat_df['customer_type'].map({'Transient': 0,␣
 ↪'Contract': 1, 'Transient-Party': 2, 'Group': 3})

cat_df['arrival_date_month']= cat_df['arrival_date_month'].map({'July' :7,␣
 ↪'August':8, 'September':9, 'October':10, 'November':11, 'December':12,␣
 ↪'January':1,
 'February':2, 'March':3, 'April':4, 'May':5, 'June':6})
```

### 2.6.2 Numerical Columns

```
[ ]: num_df = df.drop(columns = cat_cols, axis = 1)
     num_df.drop('is_canceled', axis = 1, inplace = True)
```

```
[ ]: num_df.var()
```

Based on the varaince between observations for each attibute there are some with hig variance in which we should normalize

```
[ ]: from sklearn.preprocessing import MinMaxScaler

     # List of columns to normalize
     columns_to_normalize = ['days_in_waiting_list','lead_time', 'adr', 'agent',␣
      ↪'company',␣
      ↪'arrival_date_week_number','arrival_date_day_of_month','stays_in_week_nights','previous_book

     # Initialize the MinMaxScaler with the desired range
     scaler = MinMaxScaler(feature_range=(-1, 1))

     # Fit the scaler on the data and transform the columns
     num_df[columns_to_normalize] = scaler.
      ↪fit_transform(num_df[columns_to_normalize])

     # Fill NaN values with 0
     num_df = num_df.fillna(0)
```

## 1.4 2.7 Feature Engineering

```python
num_df['total_stays'] = num_df['stays_in_weekend_nights'] +↵
 ↪num_df['stays_in_week_nights']
num_df['kids'] = num_df['children'] + num_df['babies']
cat_df['guest_location'] = cat_df['country'].apply(lambda x: 'Local' if x ==↵
 ↪'PRT' else ('International' if x != 'Unknown' else 'Unknown'))
cat_df['guest_location']= cat_df['guest_location'].map({'Unknown' :-1,'Local' :
 ↪0, 'International':1})
```

```python
num_useless=['stays_in_weekend_nights','stays_in_week_nights','children','babies']
cat_useless=['country']

num_df.drop(num_useless, axis = 1, inplace = True)
cat_df.drop(cat_useless, axis = 1, inplace = True)
```

From most booking sites, usually there are only guest and child categories (under 17 years old), so we can combine the values of children and babies into the kids column

Because the number of country values is very large in this dataset, we will condense the values to Local or International

To determine the value based on the following criteria, because this hotel dataset is located in Portugal, customers from Portugal will be set to Local and in addition to International

## 1.5 3. Model Fitting and Tuning

```python
X=pd.concat([cat_df, num_df], axis = 1)
y =df['is_canceled']
```

```python
# splitting data into training set and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30)
```

### 1.5.1 3.1 Model Fitting

**Logistic Regression:** Logistic Regression is a simple linear model used for binary classification tasks. It works by computing the probability of an instance belonging to a certain class. It uses the logistic function (sigmoid function) to map the linear combination of input features to a probability value between 0 and 1. The model then predicts the class label based on this probability value (e.g., if the probability is greater than 0.5, it predicts the positive class, otherwise the negative class).

**Random Forest Classifier:** Random Forest is an ensemble learning method that constructs multiple decision trees and combines their predictions. This results in a more robust model with better generalization capabilities. Random Forest works by creating multiple trees, each trained on a random subset of the dataset (with replacement), and using a random subset of the features at each split. The final prediction is obtained by averaging the predictions of all the trees.

**Neural Network:** A neural network is a powerful and flexible model that can learn patterns and representations in the data. It consists of layers of interconnected neurons, which are organized in a sequential structure. Each neuron receives input from the previous layer, applies a non-linear activation function, and passes the result to the next layer. Neural networks are capable of learning complex non-linear relationships in the data, making them suitable for a wide range of tasks, including image and speech recognition, natural language processing, and many others.

In the given code, we are building a neural network using Keras, with 3 fully connected (dense) hidden layers with 128, 64, and 32 neurons, respectively. The activation function used in these layers is ReLU (Rectified Linear Unit), which is a popular choice due to its non-linear nature and computational efficiency. The output layer has a single neuron with a sigmoid activation function, which outputs a probability value between 0 and 1, suitable for binary classification tasks. The model is compiled using binary cross-entropy loss and the Adam optimizer, and is then fitted on the training data. The performance of the model is evaluated on the test data using accuracy.

### 3.1.1 Logistic Regression

```
[ ]: lr = LogisticRegression()
     lr.fit(X_train, y_train)

     y_pred_lr = lr.predict(X_test)

     acc_lr = accuracy_score(y_test, y_pred_lr)
     conf = confusion_matrix(y_test, y_pred_lr)
     clf_report = classification_report(y_test, y_pred_lr)

     print(f"Accuracy Score of Logistic Regression is : {acc_lr}")
     print(f"Confusion Matrix : \n{conf}")
     print(f"Classification Report : \n{clf_report}")
```

### 3.1.2 Random Forest

```
[ ]: dtc = RandomForestClassifier()
     dtc.fit(X_train, y_train)

     y_pred_dtc = dtc.predict(X_test)

     acc_dtc = accuracy_score(y_test, y_pred_dtc)
     conf = confusion_matrix(y_test, y_pred_dtc)
     clf_report = classification_report(y_test, y_pred_dtc)

     print(f"Accuracy Score of Decision Tree is : {acc_dtc}")
     print(f"Confusion Matrix : \n{conf}")
     print(f"Classification Report : \n{clf_report}")
```

### 3.1.3 Deep Neural Network

```
[ ]: X_train_n, X_test_n, y_train_n, y_test_n = train_test_split(X, y, test_size = 0.
     ↪30)
```

```
[ ]: import keras
     from keras.layers import Dense, Dropout, BatchNormalization
     from keras.models import Sequential
     from keras.optimizers import Adam




     # Build the neural network
     model = Sequential()
     model.add(Dense(128, input_dim=X_train_n.shape[1], activation='relu'))
     model.add(Dense(64, activation='relu'))
     model.add(Dense(32, activation='relu'))
     model.add(Dense(1, activation='sigmoid'))

     # Compile the model
     optimizer = Adam(learning_rate=0.001)
     model.compile(loss='binary_crossentropy', optimizer=optimizer,␣
      ↪metrics=['accuracy'])

     model_history = model.fit(X_train_n, y_train_n, validation_data=(X_test_n,␣
      ↪y_test_n),
                               epochs=100, batch_size=32)
```

```
[99]: plt.figure(figsize = (12, 6))

      train_loss = model_history.history['loss']
      val_loss = model_history.history['val_loss']
      epoch = range(1, 101)

      loss = pd.DataFrame({'train_loss' : train_loss, 'val_loss' : val_loss})

      #px.line(data_frame = loss, x = epoch, y = ['val_loss', 'train_loss'], title =␣
       ↪'Training and Validation Loss',
              # template = 'plotly_dark')
      display.Image("newplot (4).png")
```
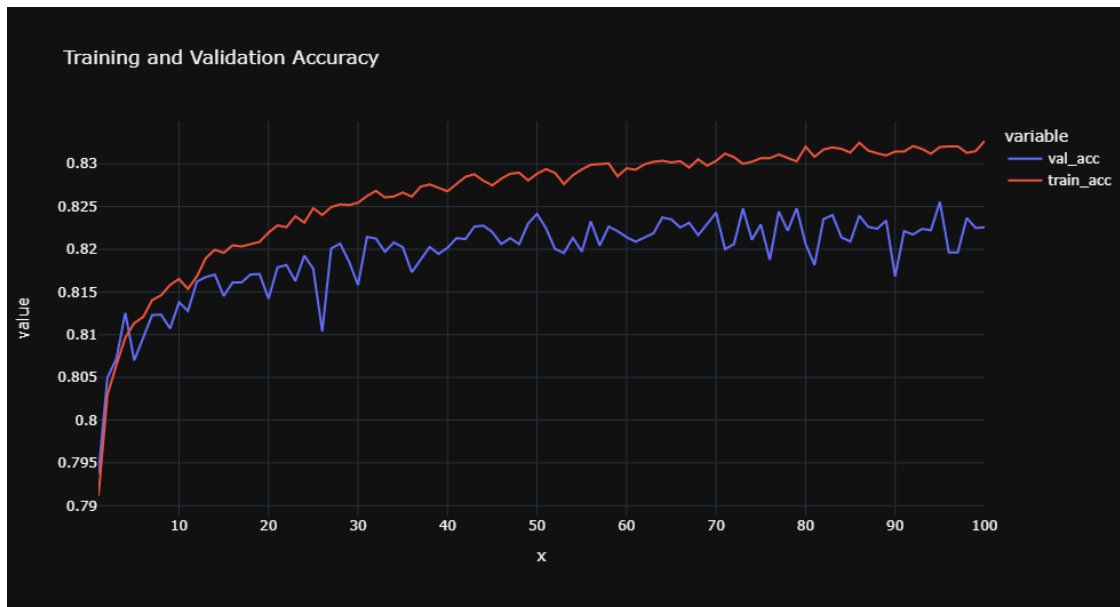
[99]:

Training and Validation Loss

```
<Figure size 960x480 with 0 Axes>
```

[100]: 
```
plt.figure(figsize = (12, 6))

train_acc = model_history.history['accuracy']
val_acc = model_history.history['val_accuracy']
epoch = range(1, 101)


accuracy = pd.DataFrame({'train_acc' : train_acc, 'val_acc' : val_acc})

#px.line(data_frame = accuracy, x = epoch, y = ['val_acc', 'train_acc'], title
 ↪= 'Training and Validation Accuracy',
     # template = 'plotly_dark')
display.Image("newplot (5).png")
```

[100]:

Training and Validation Accuracy

```
<Figure size 960x480 with 0 Axes>
```

```
[ ]: acc_ann = model.evaluate(X_test_n, y_test_n)[1]

     print(f'Accuracy of model is {acc_ann}')
```

### 1.5.2 3.2 Model Tuning

### 3.2. Logistic Regression

```
[ ]: from sklearn.model_selection import GridSearchCV

     # Define hyperparameters for tuning
     param_grid_lr = {'C': [ 0.0001, 0.001,0.01,0.1,1,10],
                      'penalty': ['l1', 'l2']}

     # Grid search
     lr = LogisticRegression()
     grid_search_lr = GridSearchCV(lr, param_grid_lr, cv=5, scoring='accuracy')
     grid_search_lr.fit(X_train, y_train)

     # Get the best hyperparameters
     best_params_lr = grid_search_lr.best_params_
     print("Best parameters for Logistic Regression:", best_params_lr)
```

```
[ ]: y_pred_lr = grid_search_lr.predict(X_test)
     acc_lr = accuracy_score(y_test, y_pred_lr)
     conf = confusion_matrix(y_test, y_pred_lr)
```

22

```python
clf_report = classification_report(y_test, y_pred_lr)

print(f"Accuracy Score of Logistic Regression is : {acc_lr}")
print(f"Confusion Matrix : \n{conf}")
print(f"Classification Report : \n{clf_report}")
```

### 3.2.2 Ramdm Forest Tree

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV

# Define hyperparameters for tuning
param_grid_rf = {'n_estimators': [10, 50, 100],
                 'criterion': ['gini'],
                 'max_depth': [1,2,5,10],
                 'min_samples_split': [1,2,5,10],
                 'min_samples_leaf': [2, 5,10]}

# Grid search
rf = RandomForestClassifier()
grid_search_rf = GridSearchCV(rf, param_grid_rf, cv=5, scoring='accuracy')
grid_search_rf.fit(X_train, y_train)

# Get the best hyperparameters
best_params_rf = grid_search_rf.best_params_
print("Best parameters for Random Forest Classifier:", best_params_rf)
```

```python
y_pred_dc = grid_search_rf.predict(X_test)
acc_dc = accuracy_score(y_test, y_pred_dc)
conf = confusion_matrix(y_test, y_pred_dc)
clf_report = classification_report(y_test, y_pred_dc)

print(f"Accuracy Score of Logistic Regression is : {acc_dc}")
print(f"Confusion Matrix : \n{conf}")
print(f"Classification Report : \n{clf_report}")
```

### 3.2.3 Deep Neural Network

```python
import keras
from keras.models import Sequential
from keras.layers import Dense
from keras.wrappers.scikit_learn import KerasClassifier
from sklearn.model_selection import GridSearchCV

X_train_n, X_test_n, y_train_n, y_test_n = train_test_split(X, y, test_size = 0.
 ↪30)
```

```python
# Function to create the Keras model
def create_model(optimizer='adam', activation='relu', neurons=100):
    model = Sequential()
    model.add(Dense(neurons, activation=activation, input_shape=(24,)))
    model.add(Dense(neurons, activation=activation))
    model.add(Dense(1, activation='sigmoid'))
    model.compile(optimizer=optimizer, loss='binary_crossentropy',␣
 ↪metrics=['accuracy'])
    return model

# Define hyperparameters for tuning
param_grid_nn = {'optimizer': ['adam',],
                 'activation': ['relu'],
                 'neurons': [50],
                 'batch_size': [64],
                 'epochs': [100]}

# Grid search
keras_classifier = KerasClassifier(build_fn=create_model, verbose=0)
grid_search_nn = GridSearchCV(keras_classifier, param_grid_nn, cv=5,␣
 ↪scoring='accuracy')
grid_search_nn.fit(X_train_n, y_train_n)

# Get the best hyperparameters
best_params_nn = grid_search_nn.best_params_
print("Best parameters for Keras Neural Network:", best_params_nn)
```

```python
from sklearn.metrics import accuracy_score, confusion_matrix,␣
 ↪classification_report

# Make predictions on the test set
y_pred_nn = grid_search_nn.predict(X_test_n)

# Calculate accuracy
acc_nn = accuracy_score(y_test_n, y_pred_nn)

# Calculate confusion matrix
conf_nn = confusion_matrix(y_test_n, y_pred_nn)

# Calculate classification report
clf_report_nn = classification_report(y_test_n, y_pred_nn)

print(f"Accuracy Score of Keras Neural Network is : {acc_nn}")
print(f"Confusion Matrix : \n{conf_nn}")
print(f"Classification Report : \n{clf_report_nn}")
```

### 1.5.3 3.3 Feature Importance

```
[82]: # Compute feature importances
      importances = dtc.feature_importances_
      indices = np.argsort(importances)[::-1]

      # Print the feature ranking
      print("Feature ranking:")
      for f in range(X.shape[1]):
          print(f"{f + 1}. {X.columns[indices[f]]} ({importances[indices[f]]:.4f})")

      # Plot the feature importances
      import matplotlib.pyplot as plt

      plt.figure()
      plt.title("Feature importances")
      plt.bar(range(X.shape[1]), importances[indices], color="r", align="center")
      plt.xticks(range(X.shape[1]), X.columns[indices], rotation=90)
      plt.xlim([-1, X.shape[1]])
      plt.show()
```
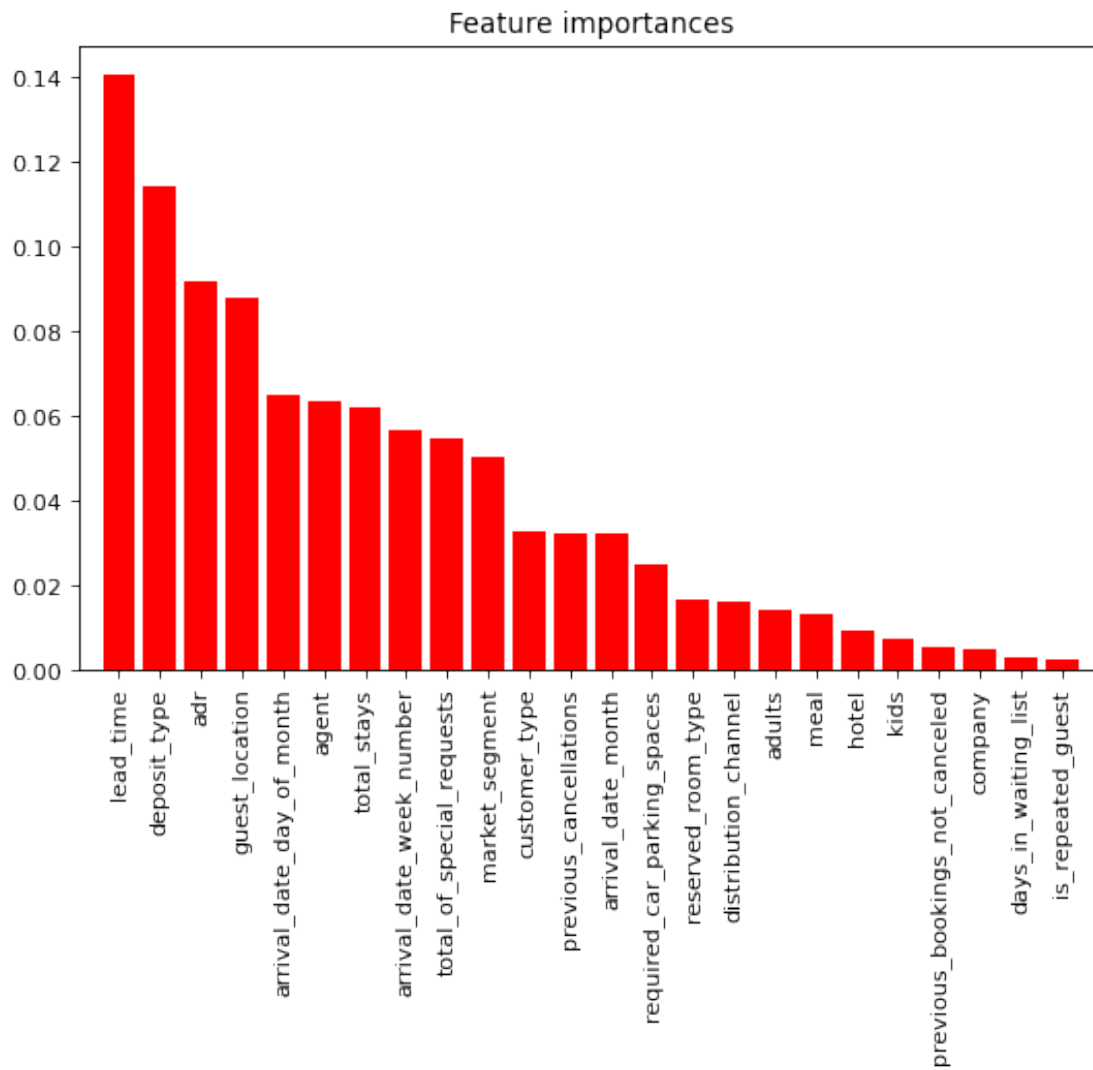
```
Feature ranking:
1. lead_time (0.1406)
2. deposit_type (0.1145)
3. adr (0.0920)
4. guest_location (0.0877)
5. arrival_date_day_of_month (0.0650)
6. agent (0.0635)
7. total_stays (0.0622)
8. arrival_date_week_number (0.0568)
9. total_of_special_requests (0.0547)
10. market_segment (0.0501)
11. customer_type (0.0327)
12. previous_cancellations (0.0323)
13. arrival_date_month (0.0321)
14. required_car_parking_spaces (0.0246)
15. reserved_room_type (0.0165)
16. distribution_channel (0.0161)
17. adults (0.0143)
18. meal (0.0131)
19. hotel (0.0091)
20. kids (0.0071)
21. previous_bookings_not_canceled (0.0051)
22. company (0.0047)
23. days_in_waiting_list (0.0027)
24. is_repeated_guest (0.0023)
```

## Feature importances



```
[ ]:  #List of the top 10 most important features
      top_features = [
          'lead_time',
          'deposit_type',
          'adr',
          'guest_location',
          'arrival_date_day_of_month',
          'agent',
          'total_stays',
          'arrival_date_week_number',
          'total_of_special_requests',
          'market_segment'
      ]
```

```
# Subset the data to include only the top 10 most important features
X_top_features = X[top_features]
```

### 3.3.1 Logistic Regression with Feature Importance

```
[ ]: X_train, X_test, y_train, y_test = train_test_split(X_top_features, y,␣
     ↪test_size = 0.30)

     lr = LogisticRegression()
     lr.fit(X_train, y_train)

     y_pred_lr = lr.predict(X_test)

     acc_lr = accuracy_score(y_test, y_pred_lr)
     conf = confusion_matrix(y_test, y_pred_lr)
     clf_report = classification_report(y_test, y_pred_lr)

     print(f"Accuracy Score of Logistic Regression is : {acc_lr}")
     print(f"Confusion Matrix : \n{conf}")
     print(f"Classification Report : \n{clf_report}")
```

### 3.3.2 DNN With Feature Importance

```
[ ]: import keras
     from keras.layers import Dense, Dropout, BatchNormalization
     from keras.models import Sequential
     from keras.optimizers import Adam

     X_train_n, X_test_n, y_train_n, y_test_n = train_test_split(X_top_features, y,␣
     ↪test_size = 0.30)


     # Build the neural network
     model = Sequential()
     model.add(Dense(128, input_dim=X_train_n.shape[1], activation='relu'))
     model.add(Dense(64, activation='relu'))
     model.add(Dense(32, activation='relu'))
     model.add(Dense(1, activation='sigmoid'))

     # Compile the model
     optimizer = Adam(learning_rate=0.005)
     model.compile(loss='binary_crossentropy', optimizer=optimizer,␣
     ↪metrics=['accuracy'])

     model_history = model.fit(X_train_n, y_train_n, validation_data=(X_test_n,␣
     ↪y_test_n),
                               epochs=100, batch_size=32)
```

27

The 100th epoch results are poch 100/100 2594/2594 [==============================]
- 6s 2ms/step - loss: 0.3563 - accuracy: 0.8327 - val_loss: 0.3822 - val_accuracy: 0.8226

## 1.6   4. Discussion & Conclusions

Based on our expirment the best perfroming model is the random forest with training set accuracy of 0.88 and testing accuracy of 0.86.

We dentiy that random forest (default parameter) as the best performing model

This model classifies whether the order will be canceled or not with an accuracy of up to 86%

As a result, this model allows hotels to estimate their occupancy rates more measurably, manage their business appropriately, and further increase their revenue.

The risks of this model are as follows:

6.4% chance that the hotel does not suspect guests, there is a risk of overbooking. 9.3% chance of a hotel allocating their resources to a reservation that will go wrong (cancelled booking) Again, the hotel also needs to be responsive and proactive in anticipating this possibility.

According to the EDA, several factors affect whether a hotel booking is canceled or not:

1-Lead Time *The longer the lead time, the higher the cancellation rate

2-Deposit Type *Cancellation rate for non-refundable deposit types reaches 2%

3-Previous Cancellations *Customers who have canceled previous hotel bookings are 91% likely to cancel again

4-Market Segment *Cancellation rate in the market segment group reaches 58%

5-Guest Location Local *Cancellation rate for local customers reaches 53%

6-Required Car Parking Space *Customers who ask for the need for a car park tend not to cancel the hotel booking

7-Special Request *Customers who ask for a special request are 95% likely to not cancel the hotel booking

8-Repeated Guest *Repeated guest 85% chance of not canceling hotel booking

9-Booking Changes *Customers who make changes to booking details are 84% likely to not cancel the hotel booking

Based on Feature Importance, top 5 categroical factors that affect hotel bookings being canceled or not:

```
1-Deposit Type - Non Refund
2-Market Segment - Online TA
3-Required Car Parking Space
4-Guest Location/Origin - International
5-Previous Cancellations
```

Business Recommendations:

```
1-For hotel bookings with a lead time of more than 7 months,
it is necessary to apply a non-refundable deposit type

2-The hotel can use the reference 54 days before the customer's
arrival date as the first step in interacting with customers who
are predicted by the model to cancel the booking (for hotel booking lead time > 54 days)

3-Implement non-refundable deposit types or advance payments
for group bookings or customers who have a history of cancellations

4-Increase hotel booking transactions through direct
distribution channels. Can be given special offer
```

## 1.7   5 Refrences

Antonio, N., Almeida, A., & Nunes, L. (2019). Predicting hotel booking cancellations using machine learning techniques. International Journal of Hospitality Management, 82, 156-164.

Schwartz, Z., Uğur, N. G., Chen, C. M., & Almeras, A. (2016). Hotel cancellations and no-shows: An application of machine learning techniques. Tourism Economics, 22(6), 1161-1177.

Tajeddini, K., & Trueman, M. (2020). Factors influencing booking cancellations in luxury hotels. International Journal of Contemporary Hospitality Management, 32(4), 1363-1383.

Zakhary, A., Atiya, A. F., & El-Shishiny, H. (2011). Forecasting hotel cancellation rates using machine learning techniques. Journal of Revenue and Pricing Management, 10(4), 344-359.

[ ]: