

MACHINE LEARNING

TOPIC – Analyzing Hotel Ratings Using Machine Learning

ABSTRACT

This research aims to develop a machine learning model to predict hotel listing discounts based on various dependent variables such as location, rating, and price. The dataset used in this study was obtained by scraping hotel listings from the official website of Trivago using the Selenium web automation tool. After preprocessing the data, a regression model was trained to predict the discount percentage for each hotel listing. The model's performance was evaluated using various regression metrics, including Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). Results indicate that the regression model achieves good predictive performance, suggesting the potential for machine learning to assist both travelers in identifying cost-effective accommodation options and hotel owners in optimizing pricing strategies.

INTRODUCTION

The hospitality industry, particularly the hotel sector, is highly competitive, with establishments constantly striving to attract guests and maintain profitability. One of the key strategies employed by hotels to achieve these objectives is offering discounts and promotions. These discounts can vary based on factors such as location, hotel rating, time of booking, and prevailing market conditions.

Understanding the dynamics of hotel discounts is crucial for both travelers and hotel operators. Travelers seek to find the best deals and maximize the value of their accommodation, while hotel operators aim to optimize occupancy rates and revenue generation. Predicting hotel discounts accurately can provide valuable insights for both parties, enabling travelers to make informed booking decisions and assisting hotel operators in devising effective pricing strategies.

In this study, we aim to develop a machine learning model that predicts hotel listing discounts based on various dependent variables. These variables include factors such as location, hotel rating, price, and potentially others that may influence discount levels. By leveraging machine learning techniques and analyzing historical data scraped from the official website of Trivago, we seek to uncover patterns and relationships that can help us predict hotel discounts with accuracy.

PROBLEM STATEMENT

The primary objective of this research is to develop a machine learning model capable of predicting hotel listing discounts based on various dependent variables, including location, rating, and price. The predictive model aims to provide insights into the factors influencing the level of discount offered by hotels, thereby assisting both travelers and hotel owners in making informed decisions.

Specifically, the problem statement can be summarized as follows:

Given a dataset of hotel listings scraped from the Trivago website, the task is to build a regression model that can predict the percentage of discount offered for each hotel room based on the following independent variables:

- Location: The city or region where the hotel is located.
- Rating: The average rating of the hotel as assigned by users.
- Property Type: The differentiation between a hotel and a resort.
- Price: The listed price of the hotel room.
- Discount: Target Variable which is to be predicted.

The model should be able to generalize well to unseen data and accurately predict the discount percentage for new hotel listings. The ultimate goal is to provide travelers with insights into cost-effective accommodation options and assist hotel owners in optimizing pricing strategies to attract more customers and maximize revenue.

```
In [30]: df['Discount'] = ((df['Price'] - df['Best Prices']) / df['Price']) * 100
```

DATA COLLECTION

The first step in our research process involved collecting relevant data that would serve as the foundation for building our predictive model. We obtained our dataset by scraping hotel listings from the official website of Trivago, a popular platform for comparing hotel prices and booking accommodations.

To automate the process of data collection, we utilized the Selenium web automation tool, which allowed us to programmatically interact with the Trivago website. By specifying search criteria such as location, dates, and other relevant parameters, we were able to retrieve a comprehensive dataset of hotel listings that met our criteria.

The data collected included various attributes for each hotel listing, such as:

1. Property Name: Name of the Hotel
2. Location: City or region where the hotel is situated.
3. Property Type: To understand if it is a hotel or resort.
4. Price: The listed price of the hotel room.
5. Lowest Price: The lowest price at which the hotel was available.
6. Rating: The average rating assigned to the hotel by users.
7. Amenities and Additional Info: Facilities provided.
8. Discount (target variable): The percentage of discount offered for the hotel room – which was calculated.

```
In [169]: from selenium import webdriver
```

```
In [170]: driver = webdriver.Chrome()
```

```
In [171]: from selenium.webdriver.common.by import By
```

```
In [172]: from selenium.webdriver.common.keys import Keys
```

```
In [173]: driver.get("https://www.trivago.in/en-IN/lm/hotels-go-a-india?search=200-64932;dr-20240502-20240509;rc-1-1")
```

```
In [174]: listings = driver.find_elements(By.CSS_SELECTOR, 'li[data-testid="accommodation-list-element"]')
len(listings)
```

```
Out[174]: 35
```

```
In [175]: name = []
for item in listings:
    property_name = item.find_element(By.CSS_SELECTOR, 'span[itemprop="name"]').text
    name.append(property_name)
    print(property_name)
```

```
ibis Styles Goa Calangute
ACRON CANDOLIM REGINA
Novotel Goa Resort and Spa
The Byke Old Anchor Beach Resort
Fairfield By Marriott Goa Anjuna
Ramada By Wyndham Goa Vagator
Royal Orchid Beach Resort & Spa
Azaya Beach Resort Goa
The Baga Beach Resort
Heritage Village Resort & Spa Goa
Hard Rock Hotel Goa
Doubletree By Hilton Goa - Panaji
Whispering Palms Beach
The Sequeira Goa
Country Inn & Suites By Carlson
Beleza By The Beach
Planet Hollywood Goa Beach Resort
Hilton Goa Resort
Le Méridien Goa, Calangute
Mercure Goa Devaaya Retreat
Novotel Goa Candolim Hotel
The Crown Goa
O Hotel Goa, Goa
Resort Terra Paraiso
Hilton Goa Resort Goa - At The Waterfront
```

```
In [177]: from selenium.common.exceptions import NoSuchElementException

property_types = []

for item in listings:
    try:
        property_type_element = item.find_element(By.CSS_SELECTOR, 'span[class="AccommodationType_hotelClass__01Wnl
        property_type = property_type_element.text
        print(property_type)
        property_types.append(property_type)
    except NoSuchElementException:
        print("Property type not found for this item.")

# Now property_types list contains scraped property types
```

```
Hotel
Hotel
Hotel
Property type not found for this item.
Hotel
Property type not found for this item.
Resort
Resort
Property type not found for this item.
Hotel
Hotel
Hotel
Resort
Property type not found for this item.
Hotel
Resort
Property type not found for this item.
Hotel
Hotel
Hotel
Hotel
..
..
```

```

5]: from selenium.common.exceptions import NoSuchElementException

best_prices_list = []

for item in listings:
    try:
        best_prices_element = item.find_element(By.CSS_SELECTOR, 'div[class="OtherDealsSection_section_Ilg5D Other
        best_prices = best_prices_element.text
        best_prices_list.append(best_prices)
    except NoSuchElementException:
        best_prices_list.append("Best prices not found for this item.")

# Now 'best_prices_list' contains scraped best prices for each item

```

```

6]: best_prices_list

```

```

6]: ['Our lowest price:\n₹4,910\nper night on Prestigia',
     'More prices',
     'Our lowest price:\n₹7,711\nper night on Agoda',
     'Our lowest price:\n₹1,875\nper night on Goibibo.com',
     'Our lowest price:\n₹4,649\nper night on ZenHotels.com',
     'More prices',
     'Our lowest price:\n₹6,496\nper night on ZenHotels.com',
     'Our lowest price:\n₹10,631\nper night on Goibibo.com',
     'More prices',
     'More prices',
     'Our lowest price:\n₹9,642\nper night on Goibibo.com',
     'Our lowest price:\n₹8,291\nper night on Agoda',
     'Our lowest price:\n₹6,751\nper night on Agoda',
     'More prices',
     '7.8',
     '']

```

```

3]: from selenium.common.exceptions import NoSuchElementException

prices = []

for item in listings:
    try:
        price_element = item.find_element(By.CSS_SELECTOR, 'span[class="Price_price_gzSve Price_large_cm2EH Price
        price = price_element.text
        print(price)
        prices.append(price)
    except NoSuchElementException:
        print("Price not found for this item.")

# Now 'prices' list contains scraped prices

```

```

₹4,912
₹4,463
₹8,159
₹2,455
₹4,650
₹7,279
₹6,827
₹11,883
₹10,204
₹7,863
₹9,662
₹9,422
₹6,960
₹1,360
₹5,840
₹8.465

```



```
1]: from selenium.common.exceptions import NoSuchElementException

aggregate_ratings = []

for item in listings:
    try:
        aggregate_rating_element = item.find_element(By.CSS_SELECTOR, 'span[data-testid="aggregate-rating"]')
        aggregate_rating = aggregate_rating_element.text
        print(aggregate_rating)
        aggregate_ratings.append(aggregate_rating)
    except NoSuchElementException:
        print("Aggregate rating not found for this item.")

# Now aggregate_ratings list contains scraped aggregate ratings
```

8.1 - Very good (681 reviews)
7.8 - Good (925 reviews)
8.2 - Very good (414 reviews)
5.8 (243 reviews)
8.0 - Very good (490 reviews)
8.5 - Excellent (33 reviews)
6.8 (942 reviews)
7.8 - Good (148 reviews)
7.6 - Good (68 reviews)
8.6 - Excellent (782 reviews)
8.1 - Very good (1309 reviews)
8.8 - Excellent (155 reviews)
7.7 - Good (1001 reviews)
6.3 (26 reviews)
8.0 - Very good (668 reviews)
8.6 - Excellent (994 reviews)
8.2 - Very good (322 reviews)
7.8 - Good (39 reviews)

```
189]: from selenium.common.exceptions import NoSuchElementException

additional_info_list = []

for item in listings:
    try:
        additional_info_element = item.find_element(By.CSS_SELECTOR, 'div[class="HotelHighlightsSection_highlights')
        additional_info_text = additional_info_element.text
        additional_info_list.append(additional_info_text)
    except NoSuchElementException:
        additional_info_list.append("Additional info not found for this item.")

# Now 'additional_info_list' contains scraped additional information for each item
```

```
190]: additional info list
```

```
190]: ['Friendly Staff, Convenient Amenities ',
      'Tasty Breakfast, Modern Amenities ',
      'Recreation And Relaxation, Tasty Food ',
      'Additional info not found for this item.',
      'Fitness Facilities, Good Service ',
      'Additional info not found for this item.',
      'Additional info not found for this item.',
      'Family-Friendly, Celebration Destination ',
      'Beach Paradise, Outdoor Activities ',
      'Additional info not found for this item.',
      'Comfortable Rooms, Dining Options ',
      'Leisure Facilities, Culinary Delights ',
      'Great Location, Local Exploration ',
      'Additional info not found for this item.',
```



```
[195]: nextpages = driver.find_elements(By.CSS_SELECTOR, 'button[class="NavigationBar_button__FY0s4"]')
nextpages[2].send_keys(Keys.ENTER)
```

```
[196]: listings = driver.find_elements(By.CSS_SELECTOR, 'li[data-testid="accommodation-list-element"]')
len(listings)
```

```
[196]: 35
```

```
[197]: for item in listings:
    property_name = item.find_element(By.CSS_SELECTOR, 'span[itemprop="name"]').text
    name.append(property_name)
    print(property_name)
```

```
The Park Calangute Goa
Bambolim Beach Resort
Hilton Goa Resort
Acron Waterfront Resort
Silver Sands Serenity
Resorte Marinha Dourada
Fabhotel K7 Trends With Pool, Baga Beach
Ramada By Wyndham Goa Vagator
Fairfield By Marriott Goa Calangute
The Astor Goa
Ginger Goa , Madgaon
Lazy Lagoon, Baga – A Lemon Tree Resort
Varanda do Mar
Radisson Goa Candolim
Pride Sun Village Resort Goa
Hotel Colonia Santa Maria
Grand Hyatt Goa
Prainha Resort & Cottage By The Sea
Golden Tulip Goa Candolim
The ...
```

```
[255]: # Trim the aggregate ratings and property types list
aggregate_ratings = aggregate_ratings[:len(name)]
property_types = property_types[:len(name)]

# Create the dictionary and DataFrame as before
data = {
    'Property Name': name,
    'Property Type': property_types,
    'Rating': rating,
    'Aggregate Rating': aggregate_ratings,
    'Price': prices,
    'Best Prices': best_prices_list,
    'Amenities': amenities_list,
    'Additional Information': additional_info_list,
    'Location': locations_list
}

df = pd.DataFrame(data)
print(df)
```

	Property Name	Property Type	Rating	\
0	ibis Styles Goa Calangute	Hotel	8.1	
1	ACRON CANDOLIM REGINA	Hotel	7.8	
2	Novotel Goa Resort and Spa	Hotel	8.2	
3	The Byke Old Anchor Beach Resort	Hotel	5.8	
4	Fairfield By Marriott Goa Anjuna	Resort	8.0	
..	
100	Viva		8.0	
101	Swati Hotel, Arambol, Goa		7.8	
102	The Fern Residency Miramar		8.2	
103	Hotel Nanutel		8.5	
104	Granpas Inn - Bougainvillea - A Heritage Hotel		8.8	

```
In [256]: df.to_csv('hotels_dataset.csv', index=False)
```

```
In [258]: import pandas as pd
df = pd.read_csv('hotels_dataset.csv')
df
```

Out[258]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location
0	ibis Styles Goa Calangute	Hotel	8.1	8.1 - Very good (681 reviews)	₹4,912	Our lowest price: ₹4,910\ner night on Prestigia	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute
1	ACRON CANDOLIM REGINA	Hotel	7.8	7.8 - Good (925 reviews)	₹4,463	More prices	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim
2	Novotel Goa Resort and Spa	Hotel	8.2	8.2 - Very good (414 reviews)	₹8,159	Our lowest price: ₹7,711\ner night on Agoda	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim
3	The Byke Old Anchor Beach Resort	Hotel	5.8	5.8 (243 reviews)	₹2,455	Our lowest price: ₹1,875\ner night on Goibib...	Amenities not found for this item.	Additional Info not found for this item.	Cavelossim
4	Fairfield By Marriott Goa Anjuna	Resort	8.0	8.0 - Very good (490 reviews)	₹4,650	Our lowest price: ₹4,649\ner night on ZenHot...	Amenities not found for this item.	Fitness Facilities, Good Service	Anjuna
...
100	Viva	NaN	8.0	8.0 - Very good (20 reviews)	₹2,495	Our lowest price: ₹2,101\ner night on Agoda	Amenities not found for this item.	Additional info not found for this item.	Margao
101	Swati Hotel, Arambol, Goa	NaN	7.8	7.8 - Good (79 reviews)	₹746	More prices	Amenities not found for this item.	Convenient Services, Relaxing Atmosphere	Pernem
102	The Fern Residency Miramar	NaN	8.2	8.2 - Very good (275 reviews)	₹4,035	Our lowest price: ₹3,484\ner night on MakeMy...	Amenities not found for this item.	Additional info not found for this item.	Panaji
103	Hotel Nanutel	NaN	8.5	8.5 - Excellent (440 reviews)	₹4,133	Our lowest price: ₹3,807\ner night on Goibib...	Amenities not found for this item.	Additional info not found for this item.	Margao
104	Granpas Inn - Bougainvillea - A Heritage Hotel	NaN	8.8	8.8 - Excellent (77 reviews)	₹4,132	Our lowest price: ₹3,688\ner night on Agoda	Amenities not found for this item.	Additional info not found for this item.	Anjuna

105 rows × 9 columns

DATA PREPROCESSING

Before modeling, the collected data underwent preprocessing to clean and prepare it for analysis. This involved tasks such as handling missing values, encoding categorical variables, and scaling numerical features. Additionally, the target variable representing the discount percentage was transformed to ensure compatibility with regression modeling techniques.

- HANDLING MISSING VALUES

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
```

```
In [3]: dff=df.copy()
```

```
In [4]: df.isnull().sum()
```

```
Out[4]: Property Name      0
Property Type      0
Rating      0
Aggregate Rating      0
Price      1
Best Prices      31
Amenities      0
Additional Information      0
Location      0
dtype: int64
```

```
In [5]: df.isnull().mean()*100
```

```
Out[5]: Property Name      0.000000
Property Type      0.000000
Rating      0.000000
Aggregate Rating      0.000000
Price      0.952381
Best Prices      29.523810
Amenities      0.000000
Additional Information      0.000000
Location      0.000000
dtype: float64
```


In [6]: df

Out[6]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location
0	ibis Styles Goa Calangute	Hotel	8.1	Very good	4912.0	4910.0	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute
1	ACRON CANDOLIM REGINA	Hotel	7.8	Good	4463.0	NaN	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim
2	Novotel Goa Resort and Spa	Hotel	8.2	Very good	8159.0	7711.0	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim
3	The Byke Old Anchor Beach Resort	Hotel	5.8	Average	2455.0	1875.0	Amenities not found for this item.	Additional info not found for this item.	Cavelossim
4	Fairfield By Marriott Goa Anjuna	Resort	8.0	Very good	4650.0	4649.0	Amenities not found for this item.	Fitness Facilities, Good Service	Anjuna
...
100	Viva	Hotel	8.0	Very good	2495.0	2101.0	Amenities not found for this item.	Additional info not found for this item.	Margao
101	Swati Hotel, Arambol, Goa	Hotel	7.8	Good	746.0	NaN	Amenities not found for this item.	Convenient Services, Relaxing Atmosphere	Pernem
102	The Fern Residency Miramar	Hotel	8.2	Very good	4035.0	3484.0	Amenities not found for this item.	Additional info not found for this item.	Panaji
103	Hotel Nanutel	Resort	8.5	Excellent	4133.0	3807.0	Amenities not found for this item.	Additional info not found for this item.	Margao
104	Granpas Inn - Bougainvillea - A Heritage Hotel	Hotel	8.8	Excellent	4132.0	3688.0	Amenities not found for this item.	Additional info not found for this item.	Anjuna

105 rows × 9 columns

```
In [7]: l = [col for col in df.columns if df[col].isnull().mean()<0.05 and df[col].isnull().mean()>0]
l
```

Out[7]: ['Price']

```
In [8]: df.dropna(subset=l, inplace=True)
```

```
In [9]: df['Property Type'].fillna(df['Property Type'].mode()[0], inplace=True)
df
```

Out[9]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location
0	ibis Styles Goa Calangute	Hotel	8.1	Very good	4912.0	4910.0	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute
1	ACRON CANDOLIM REGINA	Hotel	7.8	Good	4463.0	NaN	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim
2	Novotel Goa Resort and Spa	Hotel	8.2	Very good	8159.0	7711.0	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim
3	The Byke Old Anchor Beach Resort	Hotel	5.8	Average	2455.0	1875.0	Amenities not found for this item.	Additional info not found for this item.	Cavelossim
4	Fairfield By Marriott Goa Anjuna	Resort	8.0	Very good	4650.0	4649.0	Amenities not found for this item.	Fitness Facilities, Good Service	Anjuna
...
100	Viva	Hotel	8.0	Very good	2495.0	2101.0	Amenities not found for this item.	Additional info not found for this item.	Margao
101	Swati Hotel, Arambol, Goa	Hotel	7.8	Good	746.0	NaN	Amenities not found for this item.	Convenient Services, Relaxing Atmosphere	Pernem
102	The Fern Residency Miramar	Hotel	8.2	Very good	4035.0	3484.0	Amenities not found for this item.	Additional info not found for this item.	Panaji
103	Hotel Nanutel	Resort	8.5	Excellent	4133.0	3807.0	Amenities not found for this item.	Additional info not found for this item.	Margao
104	Granpas Inn - Bougainvillea - A Heritage Hotel	Hotel	8.8	Excellent	4132.0	3688.0	Amenities not found for this item.	Additional info not found for this item.	Anjuna

104 rows × 9 columns

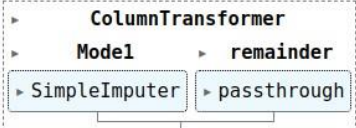
```
In [10]: imp=SimpleImputer(strategy="most_frequent")
```

```
In [11]: tr=ColumnTransformer([
          ('Model',imp,['Property Type'])],
          remainder='passthrough')
```

```
In [12]: X = df.iloc[:,0:3]
          y = df.iloc[:, -1]
```

```
In [13]: tr.fit(X)
```

```
Out[13]:
```

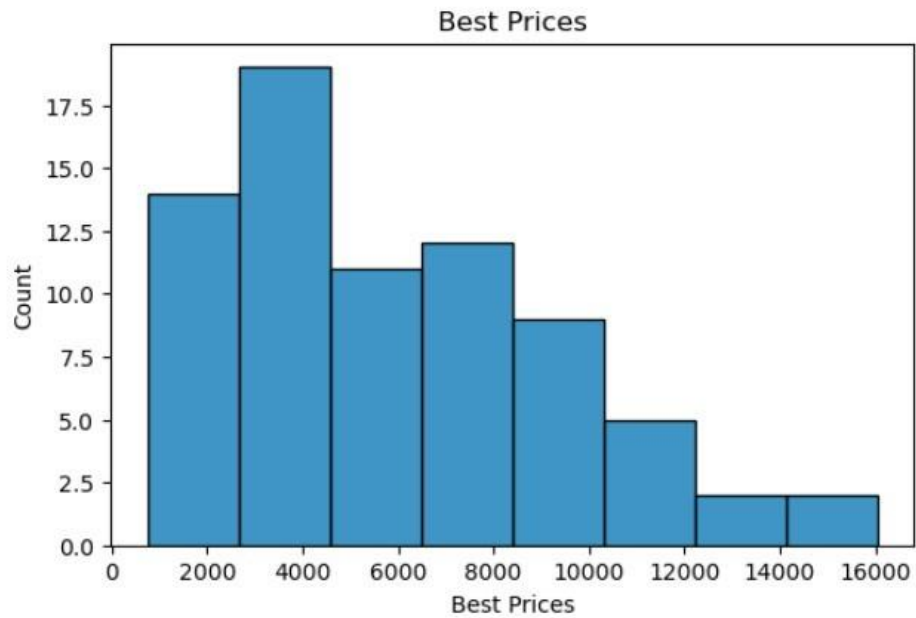


The diagram shows a dashed box labeled 'ColumnTransformer'. Inside, there are two main sections: 'Model' and 'remainder'. Under 'Model', there is a box labeled 'SimpleImputer'. Under 'remainder', there is a box labeled 'passthrough'. Arrows point from 'SimpleImputer' and 'passthrough' to a common point below them, indicating they are combined in the transformation process.

```
In [14]: X=tr.fit_transform(X)
          X
          ['Hotel', 'The Park Edge River Side', 8.0],
          ['Hotel', 'Hotel Meraden Opus', 6.6],
          ['Hotel', 'Kashinath Beach Huts', 7.8],
          ['Hotel', 'The Ivy Anjuna', 8.0],
          ['Hotel', 'Cupid's Heaven Beach Resort', 8.1],
          ['Resort', 'Nanu Resort - Arambol', 9.0],
          ['Hotel', 'Fortune Miramar', 7.5],
          ['Hotel', 'The Center Court Goa', 8.5],
          ['Hotel', 'Regenta Resort Varcas Beach Goa', 6.3],
          ['Hotel', 'Dunhill Beach Resort', 8.1],
          ['Hotel', 'De Alturas Resort', 7.1],
          ['Hotel', 'Hotel Pirache Art', 7.0],
          ['Resort', 'Spazio Leisure Resort', 8.0],
          ['Resort', 'Grandeur De Sanchi', 6.9],
          ['Hotel', 'Viva', 8.0],
          ['Hotel', 'Swati Hotel, Arambol, Goa', 7.8],
          ['Hotel', 'The Fern Residency Miramar', 8.2],
          ['Resort', 'Hotel Nanutel', 8.5],
          ['Hotel', 'Granpas Inn - Bougainvillea - A Heritage Hotel', 8.8]],
          dtype=object)
```

```
In [15]: plt.figure(figsize=(14,4))  
  
plt.subplot(121)  
sns.histplot(df['Best Prices'])  
plt.title('Best Prices')
```

Out[15]: Text(0.5, 1.0, 'Best Prices')



```
In [16]: df['Best Prices'].fillna(df['Best Prices'].mean(),inplace=True)
```



```
In [24]: df['Aggregate Rating'].fillna(df['Aggregate Rating'].mode()[0], inplace=True)
df
```

Out[24]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location
0	ibis Styles Goa Calangute	0.0	8.1	Very good	4912.0	4910.00000	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute
1	ACRON CANDOLIM REGINA	0.0	7.8	Good	4463.0	5958.72973	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim
2	Novotel Goa Resort and Spa	0.0	8.2	Very good	8159.0	7711.00000	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim
3	The Byke Old Anchor Beach Resort	0.0	5.8	Average	2455.0	1875.00000	Amenities not found for this item.	Additional info not found for this item.	Cavelossim
4	Fairfield By Marriott Goa Anjuna	1.0	8.0	Very good	4650.0	4649.00000	Amenities not found for this item.	Fitness Facilities, Good Service	Anjuna
...
100	Viva	0.0	8.0	Very good	2495.0	2101.00000	Amenities not found for this item.	Additional info not found for this item.	Margao
101	Swati Hotel, Arambol, Goa	0.0	7.8	Good	746.0	5958.72973	Amenities not found for this item.	Convenient Services, Relaxing Atmosphere	Pemem
102	The Fern Residency Miramar	0.0	8.2	Very good	4035.0	3484.00000	Amenities not found for this item.	Additional info not found for this item.	Panaji
103	Hotel Nanutel	1.0	8.5	Excellent	4133.0	3807.00000	Amenities not found for this item.	Additional info not found for this item.	Margao
104	Granpas Inn - Bougainvillea - A Heritage Hotel	0.0	8.8	Excellent	4132.0	3688.00000	Amenities not found for this item.	Additional info not found for this item.	Anjuna

104 rows × 9 columns

```
In [28]: df.drop(columns=['Amenities'], inplace=True)
```

```
In [29]: df
```

Out[29]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Additional Information	Location	Location_Frequency_Encoded
0	ibis Styles Goa Calangute	0.0	8.1	0.0	4912.0	4910.00000	Friendly Staff, Convenient Amenities	Calangute	0.163462
1	ACRON CANDOLIM REGINA	0.0	7.8	1.0	4463.0	5958.72973	Tasty Breakfast, Modern Amenities	Candolim	0.192308

```
In [32]: df.drop(columns=['Additional Information'], inplace=True)|
```

```
In [33]: df
```

Out[33]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices
0	ibis Styles Goa Calangute	0.0	8.1	0.0	4912.0	4910.00000
1	ACRON CANDOLIM REGINA	0.0	7.8	1.0	4463.0	5958.72973
2	Novotel Goa Resort and Spa	0.0	8.2	2.0	8159.0	7711.00000

- ENCODING

```
In [18]: df['Property Type'].unique()
```

```
Out[18]: array(['Hotel', 'Resort'], dtype=object)
```

```
In [19]: print("Length of DataFrame:", len(df))
print("Length of 'Property Type' column:", len(df['Property Type']))
```

```
Length of DataFrame: 104
Length of 'Property Type' column: 104
```

```
In [20]: print(df['Property Type'].value_counts())
```

```
Property Type
Hotel      79
Resort     25
Name: count, dtype: int64
```

```
In [21]: import numpy as np
from sklearn.preprocessing import OrdinalEncoder

encoder = OrdinalEncoder(categories=[['Hotel', 'Resort']])
df['Property Type'] = encoder.fit_transform(df[['Property Type']])

df['Property Type'].unique()
```

```
Out[21]: array([0., 1.])
```

```
In [25]: # Get the unique values in the 'Aggregate Rating' column and convert it to a list
categories = df['Aggregate Rating'].unique().tolist()
```

```
# Initialize OrdinalEncoder with defined categories
encoder = OrdinalEncoder(categories=[categories])

# Convert qualitative ratings to numerical categories
ratings_encoded = encoder.fit_transform(df[['Aggregate Rating']])
```

```
# Assign encoded ratings back to the DataFrame
df['Aggregate Rating'] = ratings_encoded
```

```
# Check the unique encoded values
print(df['Aggregate Rating'].unique())
```

```
[ 0.  1.  2.  3.  4.  5.  6.  7.  8.  9. 10. 11.]
```

```
In [26]: df
```

```
Out[26]:
```

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location
0	ibis Styles Goa Calangute	0.0	8.1	0.0	4912.0	4910.00000	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute
1	ACRON CANDOLIM REGINA	0.0	7.8	1.0	4463.0	5958.72973	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim
2	Novotel Goa Resort and Spa	0.0	8.2	2.0	8159.0	7711.00000	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim
3	The Byke Old Anchor Beach Resort	0.0	5.8	3.0	2455.0	1875.00000	Amenities not found for this item.	Additional info not found for this item.	Cavelossim

```
In [27]: location_frequency = df['Location'].value_counts(normalize=True)

# Map the frequency values to the corresponding categories
df['Location_Frequency_Encoded'] = df['Location'].map(location_frequency)
df
```

Out[27]:

	Property Name	Property Type	Rating	Aggregate Rating	Price	Best Prices	Amenities	Additional Information	Location	Location_Frequency_Encoded
0	Ibis Styles Goa Calangute	0.0	8.1	0.0	4912.0	4910.00000	Amenities not found for this item.	Friendly Staff, Convenient Amenities	Calangute	0.163462
1	ACRON CANDOLIM REGINA	0.0	7.8	1.0	4463.0	5958.72973	Amenities not found for this item.	Tasty Breakfast, Modern Amenities	Candolim	0.192308
2	Novotel Goa Resort and Spa	0.0	8.2	2.0	8159.0	7711.00000	Amenities not found for this item.	Recreation And Relaxation, Tasty Food	Candolim	0.192308
3	The Byke Old Anchor Beach Resort	0.0	5.8	3.0	2455.0	1875.00000	Amenities not found for this item.	Additional Info not found for this item.	Cavelossim	0.009615
4	Fairfield By Marriott Goa Anjuna	1.0	8.0	0.0	4650.0	4649.00000	Amenities not found for this item.	Fitness Facilities, Good Service	Anjuna	0.086538
...
							Amenities not	Additional info not		

• SCALING AND NORMALIZATION

```
In [34]: from sklearn.preprocessing import StandardScaler

# Initialize the StandardScaler
scaler = StandardScaler()

# Define the columns to be scaled
columns_to_scale = ['Rating', 'Aggregate Rating', 'Price', 'Best Prices', 'Location_Frequency_Encoded', 'Discount']

# Fit and transform the selected columns
df[columns_to_scale] = scaler.fit_transform(df[columns_to_scale])

# Display the scaled dataset
print(df.head())
```

	Property Name	Property Type	Rating \
0	Ibis Styles Goa Calangute	0.0	0.171563
1	ACRON CANDOLIM REGINA	0.0	-0.171563
2	Novotel Goa Resort and Spa	0.0	0.285939
3	The Byke Old Anchor Beach Resort	0.0	-2.459074
4	Fairfield By Marriott Goa Anjuna	1.0	0.057188

	Aggregate Rating	Price	Best Prices	Location \
0	-1.529665	-0.444566	-3.549830e-01	Calangute
1	-1.165626	-0.563642	-6.157071e-16	Candolim
2	-0.801587	0.416546	5.931234e-01	Candolim
3	-0.437547	-1.096168	-1.382296e+00	Cavelossim
4	-1.529665	-0.514049	-4.433285e-01	Anjuna

	Location_Frequency_Encoded	Discount
0	0.946924	0.088750
1	1.395922	-0.283748
2	1.395922	0.149253
3	-1.447729	0.350567
4	-0.250402	0.088536

MODEL BUILDING

For the regression task of predicting hotel listing discounts, various machine learning algorithms were evaluated, including linear regression, decision trees, and ensemble methods. After experimenting with different models, a Gradient Boosting Regressor was selected for its ability to capture complex relationships in the data and handle both numerical and categorical features effectively.

- KNN

```
In [38]: df.drop(columns=['Property Name'], inplace=True)
df.drop(columns=['Location'], inplace=True)
```

```
In [39]: df
```

```
Out[39]:
```

	Property Type	Rating	Aggregate Rating	Price	Best Prices	Location_Frequency_Encoded	Discount
0	0.0	0.171563	-1.529665	-0.444566	-3.549830e-01	0.946924	0.088750
1	0.0	-0.171563	-1.165626	-0.563642	-6.157071e-16	1.395922	-0.283748
2	0.0	0.285939	-0.801587	0.416546	5.931234e-01	1.395922	0.149253
3	0.0	-2.459074	-0.437547	-1.096168	-1.382296e+00	-1.447729	0.350567
4	1.0	0.057188	-1.529665	-0.514049	-4.433285e-01	-0.250402	0.088536
...
100	0.0	0.057188	-0.801587	-1.085560	-1.305797e+00	-0.400068	0.263603
101	0.0	-0.171563	0.290531	-1.549399	-6.157071e-16	-1.148398	-7.668750
102	0.0	0.285939	-0.801587	-0.677149	-8.376676e-01	0.198595	0.239890
103	1.0	0.629065	-0.073508	-0.651159	-7.283358e-01	-0.400068	0.175861
104	0.0	0.972192	-0.073508	-0.651424	-7.686160e-01	-0.250402	0.207585

104 rows × 7 columns

```
In [40]: df.shape
```

```
Out[40]: (104, 7)
```

```
In [41]: x=df.iloc[:,1:]
y=df.iloc[:,0]
```

```
In [42]: y.unique()
```

```
Out[42]: array([0., 1.])
```

```
In [43]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x,y,test_size=0.2, random_state=2)
```

```
In [44]: X_train.head()
```

Out[44]:

	Rating	Aggregate Rating	Price	Best Prices	Location_Frequency_Encoded	Discount
11	0.972192	0.654571	0.751497	7.894467e-01	0.198595	0.221555
89	0.057188	-0.801587	-0.759627	-6.157071e-16	-0.250402	-0.577873
62	-2.001572	-0.437547	-1.135153	-6.157071e-16	-0.549734	-1.667660
74	-0.743441	-0.437547	-0.960650	-1.079687e+00	1.395922	0.162031
5	0.629065	-0.073508	0.183168	-6.157071e-16	-0.250402	0.289652

```
In [45]: X_train.shape
```

Out[45]: (83, 6)

```
In [46]: from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [47]: X_train
```

Out[47]: array([[0.9203269 , 0.64093121, 0.79482526, 0.85038399, 0.24565369,
 0.2191171],
 [0.01633716, -0.83677131, -0.72511057, 0.052241 , -0.22127585,
 -0.55146207],
 [-2.01763975, -0.46734568, -1.10282716, 0.052241 , -0.53256221,
 -1.60192271],
 [-0.77465386, -0.46734568, -0.92730631, -1.03933937, 1.49079913,
 0.16174175],
 [0.58133075, -0.09792005, 0.22318214, 0.052241 , -0.22127585, ...])

```
In [48]: X_train.shape
```

```
Out[48]: (83, 6)
```

```
In [49]: from sklearn.neighbors import KNeighborsClassifier  
knn = KNeighborsClassifier(n_neighbors=3)
```

```
In [50]: knn.fit(X_train,y_train)
```

```
Out[50]: 

KNeighborsClassifier



KNeighborsClassifier(n_neighbors=3)


```

```
In [51]: from sklearn.metrics import accuracy_score  
y_pred = knn.predict(X_test)  
accuracy_score(y_test, y_pred)
```

```
Out[51]: 0.7619047619047619
```

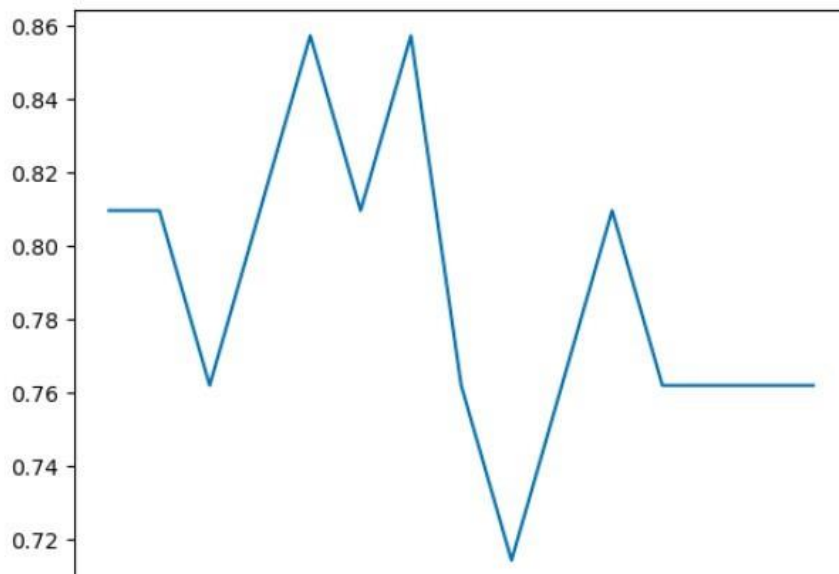
```
In [52]: y_train
```

```
Out[52]: 11    1.0  
89    0.0  
62    0.0  
74    1.0  
5     1.0  
...  
43    0.0  
22    0.0  
73    0.0  
15    0.0  
40    0.0  
Name: Property Type, Length: 83, dtype: float64
```

```
In [53]: scores = []  
  
for i in range(1,16):  
    knn = KNeighborsClassifier(n_neighbors=i)  
    knn.fit(X_train,y_train)  
    y_pred = knn.predict(X_test)  
    scores.append(accuracy_score(y_test, y_pred))
```

```
In [54]: import matplotlib.pyplot as plt  
plt.plot(range(1,16),scores)
```

```
Out[54]: [matplotlib.lines.Line2D at 0x7ee7d56c5e50]
```



```
In [55]: from sklearn import preprocessing  
  
# label_encoder object knows  
# how to understand word labels.  
label_encoder = preprocessing.LabelEncoder()  
y_trans=label_encoder.fit_transform(y_train)
```

```
In [56]: from sklearn.decomposition import PCA  
  
pca = PCA(n_components = 2)  
X_train2 = pca.fit_transform(X_train)  
knn.fit(X_train2,y_trans)
```

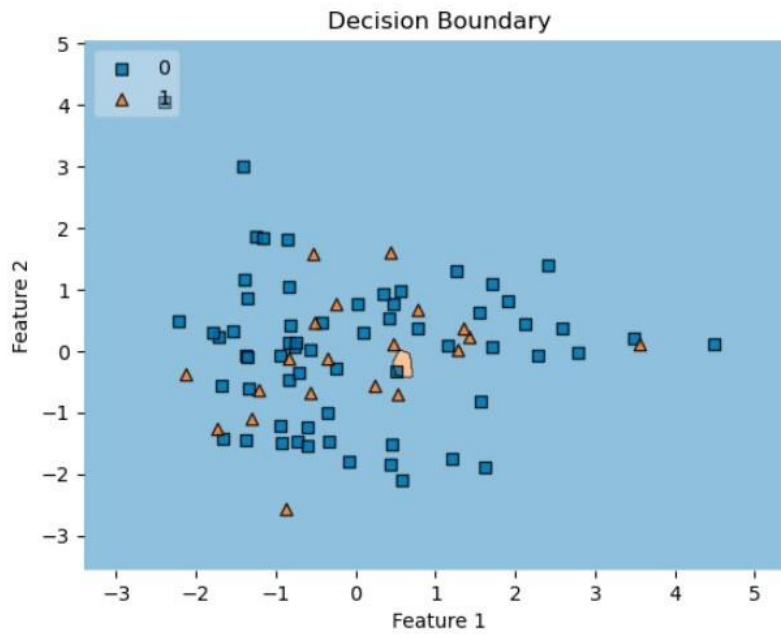
```
Out[56]: KNeighborsClassifier  
KNeighborsClassifier(n_neighbors=15)
```



```
In [58]: import matplotlib.pyplot as plt
from mlxtend.plotting import plot_decision_regions # Install mlxtend library if not installed

# Visualize decision boundary for a 2D dataset

plot_decision_regions(X_train2, y_train, clf=knn, legend=2)
plt.xlabel("Feature 1")
plt.ylabel("Feature 2")
plt.title("Decision Boundary")
plt.show()
```



- LOGISTIC REGRESSION

```
In [60]: from sklearn.model_selection import train_test_split
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import MinMaxScaler
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.feature_selection import SelectKBest, chi2
```

```
In [61]: clf = LogisticRegression()
```

```
In [62]: pipe = Pipeline([
    ('clf', clf)
])
```

```
In [63]: # train
pipe.fit(X_train, y_train)
```

```
Out[63]: Pipeline
  LogisticRegression
  LogisticRegression()
```

```
In [64]: # Display Pipeline
from sklearn import set_config
set_config(display='diagram')

In [65]: # Predict
y_pred = pipe.predict(X_test)

In [66]: from sklearn.metrics import accuracy_score
accuracy_score(y_test, y_pred)

Out[66]: 0.7619047619047619
```

Cross Validation

```
In [67]: # cross validation using cross_val_score
from sklearn.model_selection import cross_val_score
cross_val_score(pipe, X_train, y_train, cv=5, scoring='accuracy').mean()

Out[67]: 0.7470588235294118
```

Grid Search CV

```
In [68]: # gridsearchcv
param_grid = {
    "clf_penalty": ['l1', 'l2'],
    "clf_C": [0.001, 0.01, 0.1, 1, 10, 100]
}
```

```
In [69]: from sklearn.model_selection import GridSearchCV
grid = GridSearchCV(pipe, param_grid, cv=5, scoring='accuracy')
grid.fit(X_train, y_train)
```

```
Out[69]:
GridSearchCV
└─ estimator: Pipeline
   └─ LogisticRegression
```

```
In [70]: grid.best_score_

Out[70]: 0.7588235294117647
```

```
In [71]: grid.best_params_

Out[71]: {'clf_C': 0.001, 'clf_penalty': 'l2'}
```

Export to a pickle file

```
In [72]: # export
import pickle
pickle.dump(pipe, open('lr_pile.pkl', 'wb'))
```

- RANDOM FOREST

```
In [76]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression

from sklearn.metrics import accuracy_score

In [77]: X = df.iloc[:,0:-1]
y = df.iloc[:, -1]

In [78]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

In [80]: rf = RandomForestRegressor(oob_score=True)
rf.fit(X_train, y_train)

# Get the out-of-bag score
rf.oob_score_

Out[80]: 0.44407300162749097

In [82]: from sklearn.metrics import mean_squared_error

mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error:", mse)

Mean Squared Error: 0.016720605598622634
```

Testing with multiple models

```
In [83]: rf = RandomForestClassifier()
svc = SVC()
lr = LogisticRegression()

In [85]: # Initialize and fit the RandomForestRegressor model
rf_regressor = RandomForestRegressor()
rf_regressor.fit(X_train, y_train)

# Predict the target variable for the test data
y_pred = rf_regressor.predict(X_test)

# Evaluate the performance of the regression model using appropriate metrics
# For example, you can calculate the mean squared error (MSE)
from sklearn.metrics import mean_squared_error
mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error:", mse)

Mean Squared Error: 0.014002850060311914

In [92]: from sklearn.ensemble import RandomForestRegressor

# Initialize and fit the RandomForestRegressor model
rf = RandomForestRegressor(max_samples=0.75, random_state=42)
rf.fit(X_train, y_train)

# Predict the target variable for the test data
y_pred = rf.predict(X_test)
```

```
In [94]: from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestRegressor

# Perform cross-validation with RandomForestRegressor
np.mean(cross_val_score(RandomForestRegressor(max_samples=0.75), X, y, cv=10, scoring='neg_mean_squared_error'))
```

Out[94]: -0.5865336105780168

Grid Search CV

```
In [96]: # Number of trees in random forest
n_estimators = [20,60,100,120]

# Number of features to consider at every split
max_features = [0.2,0.6,1.0]

# Maximum number of levels in tree
max_depth = [2,8,None]

# Number of samples
max_samples = [0.5,0.75,1.0]

# 108 diff random forest train
```

```
In [97]: param_grid = {'n_estimators': n_estimators,
                      'max_features': max_features,
                      'max_depth': max_depth,
                      'max_samples': max_samples
                    }
print(param_grid)

{'n_estimators': [20, 60, 100, 120], 'max_features': [0.2, 0.6, 1.0], 'max_depth': [2, 8, None], 'max_samples': [0.5, 0.75, 1.0]}
```

```
In [98]: rf = RandomForestRegressor()
```

```
In [99]: from sklearn.model_selection import GridSearchCV

rf_grid = GridSearchCV(estimator = rf,
                      param_grid = param_grid,
                      cv = 5,
                      verbose=2,
                      n_jobs = -1)
```



```
In [100]: rf_grid.fit(X_train,y_train)
```

Fitting 5 folds for each of 108 candidates, totalling 540 fits

```
Out[100]:
```

```
GridSearchCV
  estimator: RandomForestRegressor
    RandomForestRegressor
```

```
In [101]: rf_grid.best_params_
```

```
Out[101]: {'max_depth': 8, 'max_features': 1.0, 'max_samples': 1.0, 'n_estimators': 100}
```

```
In [102]: rf_grid.best_score_
```

```
Out[102]: 0.5911837456533983
```

```
[CV] END max_depth=2, max_features=0.2, max_samples=0.5, n_estimators=60; total time= 0.1s
[CV] END max_depth=2, max_features=0.2, max_samples=0.5, n_estimators=100; total time= 0.2s
[CV] END max_depth=2, max_features=0.2, max_samples=0.75, n_estimators=20; total time= 0.0s
[CV] END max_depth=2, max_features=0.2, max_samples=0.75, n_estimators=60; total time= 0.1s
[CV] END max_depth=2, max_features=0.2, max_samples=0.75, n_estimators=100; total time= 0.1s
[CV] END max_depth=2, max_features=0.2, max_samples=1.0, n_estimators=20; total time= 0.0s
[CV] END max_depth=2, max_features=0.2, max_samples=1.0, n_estimators=100; total time= 0.3s
[CV] END max_depth=2, max_features=0.6, max_samples=0.5, n_estimators=20; total time= 0.0s
[CV] END max_depth=2, max_features=0.6, max_samples=0.5, n_estimators=60; total time= 0.1s
[CV] END max_depth=2, max_features=0.6, max_samples=0.5, n_estimators=100; total time= 0.2s
[CV] END max_depth=2, max_features=0.6, max_samples=0.5, n_estimators=120; total time= 0.2s
[CV] END max_depth=2, max_features=0.6, max_samples=0.75, n_estimators=120; total time= 0.2s
[CV] END max_depth=2, max_features=0.6, max_samples=1.0, n_estimators=60; total time= 0.1s
[CV] END max_depth=2, max_features=0.6, max_samples=1.0, n_estimators=120; total time= 0.2s
[CV] END max_depth=2, max_features=1.0, max_samples=0.5, n_estimators=60; total time= 0.1s
[CV] END max_depth=2, max_features=1.0, max_samples=0.5, n_estimators=120; total time= 0.2s
[CV] END max_depth=2, max_features=1.0, max_samples=0.75, n_estimators=60; total time= 0.1s
```

- NAIVE BAYES

```
In [103]: # Importing necessary libraries
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
```

```
In [105]: # Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

```
In [106]: # Initialize Gaussian Naive Bayes classifier
model = GaussianNB()

# Train the model
model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)
```

```
In [107]: # Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)

Accuracy: 0.9777777777777777
```

```
In [109]: from sklearn.naive_bayes import MultinomialNB
# Initialize Multinomial Naive Bayes classifier
model = MultinomialNB()

# Train the model
model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)
```

```
In [110]: # Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

Accuracy: 0.9555555555555556

- SVM

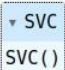
```
In [111]: import numpy as np
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score
```

```
In [112]: # Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

```
In [113]: # Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [114]: # Create SVM classifier
svm_classifier = SVC()

# Train the classifier
svm_classifier.fit(X_train, y_train)
```

Out[114]: 

```
In [115]: # Predict using the trained model
y_pred = svm_classifier.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

Accuracy: 1.0

```
In [116]: y_pred
```

Out[116]: array([1, 0, 2, 1, 1, 0, 1, 2, 1, 1, 2, 0, 0, 0, 0, 1, 2, 1, 1, 2, 0, 2,
0, 2, 2, 2, 2, 2, 0, 0, 0, 0, 1, 0, 0, 2, 1, 0, 0, 0, 2, 1, 1, 0,
0])

MODEL EVALUATION

The performance of the trained regression model was evaluated using standard regression metrics, including Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). These metrics provide insights into the accuracy and precision of the model's predictions compared to the actual discount percentages observed in the dataset. Additionally, visualizations such as scatter plots of predicted versus actual discounts were used to assess the model's performance across different ranges of discount values.

Comparing Multiple Models using AUC-ROC Curve

```
In [121]: from sklearn.linear_model import LogisticRegression
```

```
# Initialize the Logistic Regression model  
lr = LogisticRegression()  
  
# Fit the model to your training data  
lr.fit(X_train, y_train)  
  
# Now you can use the predict_proba method  
y_scores = lr.predict_proba(X_test)[:, 1]
```

```
In [126]: y_pred = model.predict(X_test)  
y_true = y_test # Assuming y_test contains the true labels corresponding to X_test
```

```
In [128]: class_names = np.unique(np.concatenate((y_true, y_pred)))
```

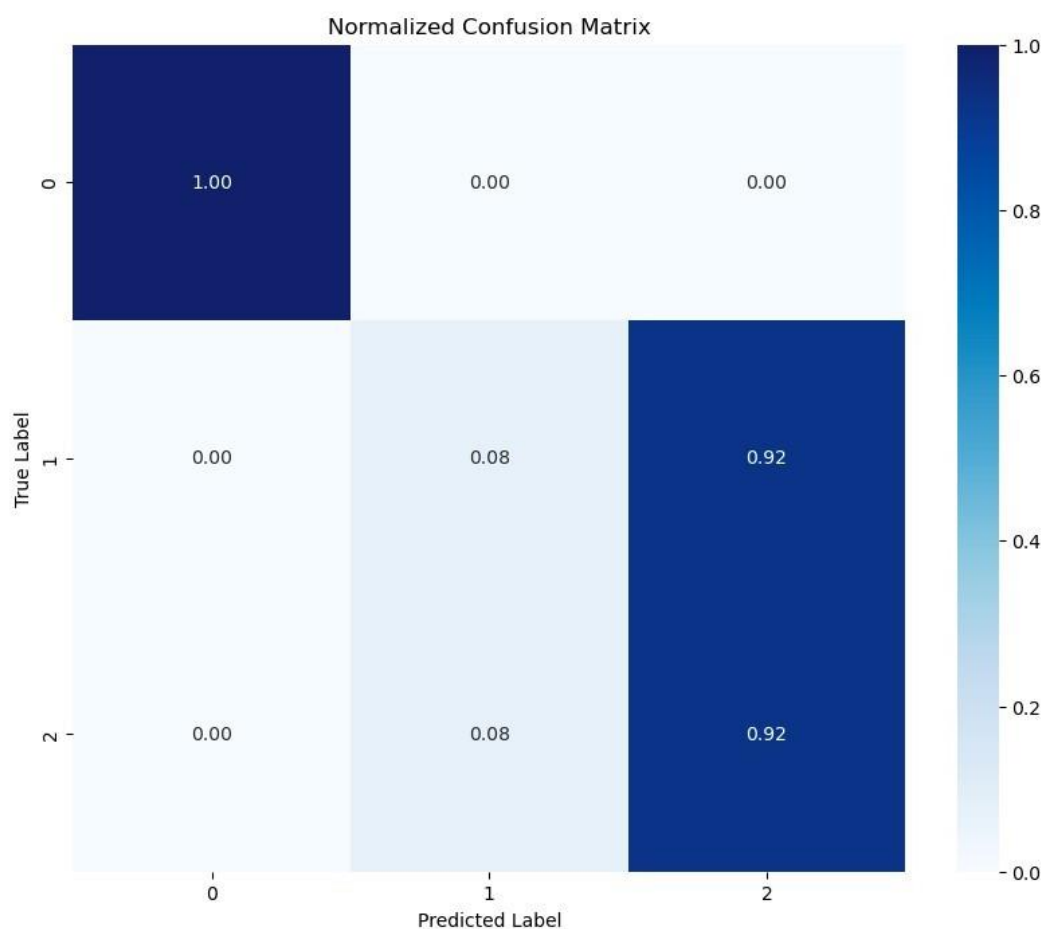


```
In [129]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix

# Assuming y_true contains true labels and y_pred contains predicted labels
# Compute confusion matrix
cm = confusion_matrix(y_true, y_pred)

# Normalize confusion matrix
cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]

# Plot confusion matrix heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(cm_normalized, annot=True, cmap='Blues', fmt='.2f',
            xticklabels=class_names, yticklabels=class_names)
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.title('Normalized Confusion Matrix')
plt.show()
```



```
In [137]: from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import LogisticRegression

# Assuming you have a logistic regression classifier 'lr' and test data 'X_test', and labels 'y_test'
# Create an instance of OneVsRestClassifier
ovr = OneVsRestClassifier(LogisticRegression())

# Fit the classifier
ovr.fit(X_train, y_train)

# Predict probabilities for each class using OneVsRestClassifier
y_scores_ovr = ovr.predict_proba(X_test)

# Now you can compute ROC curve and ROC area for each class and plot them as needed
```

```
In [141]: n_classes = len(np.unique(y_test))
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt
from sklearn.preprocessing import label_binarize
from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import LogisticRegression
import numpy as np

# Assuming you have a logistic regression classifier 'lr' and test data 'X_test', and labels 'y_test'
# Create an instance of OneVsRestClassifier
ovr = OneVsRestClassifier(LogisticRegression())

# Fit the classifier
ovr.fit(X_train, y_train)

# Predict probabilities for each class using OneVsRestClassifier
y_scores_ovr = ovr.predict_proba(X_test)

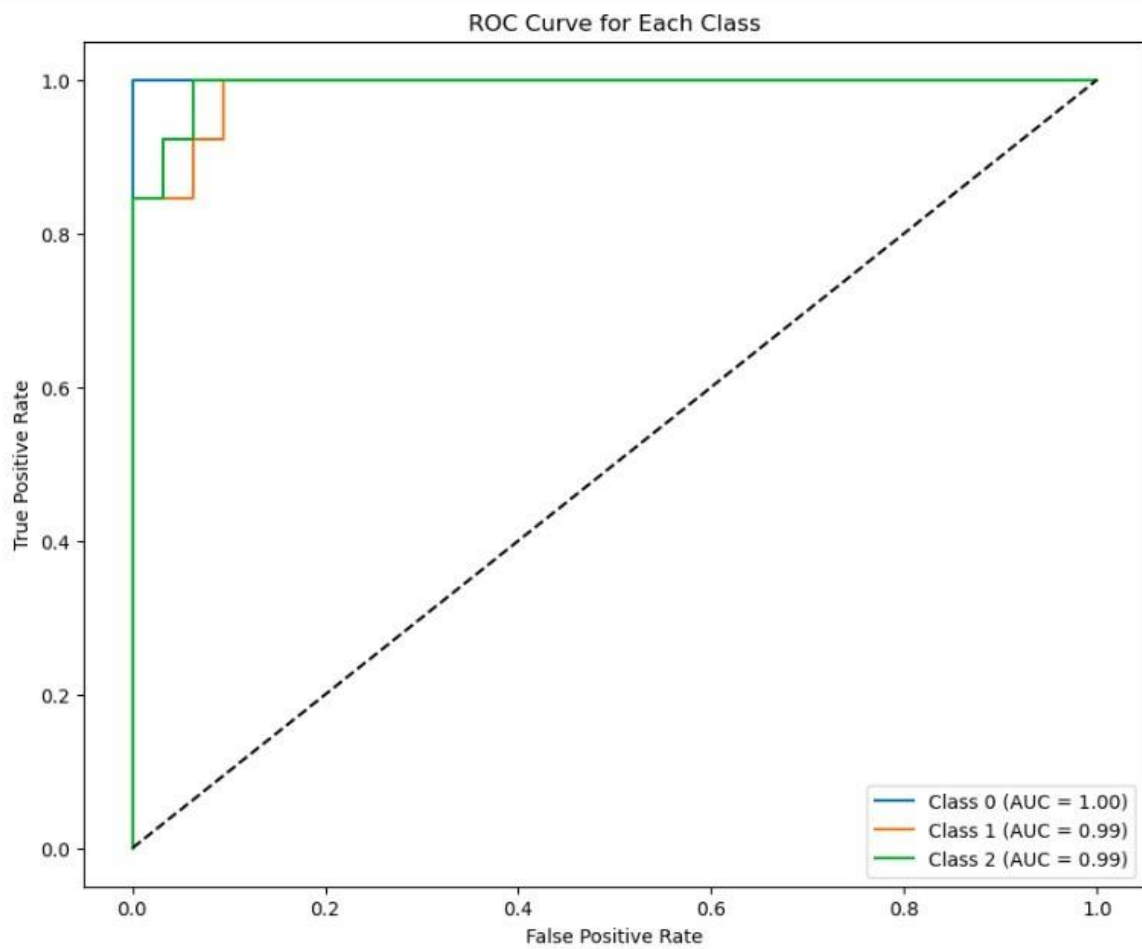
# Binarize the labels
y_test_binarized = label_binarize(y_test, classes=np.unique(y_test))

# Define the number of classes
n_classes = len(np.unique(y_test))

# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(n_classes):
    fpr[i], tpr[i], _ = roc_curve(y_test_binarized[:, i], y_scores_ovr[:, i])
    roc_auc[i] = auc(fpr[i], tpr[i])

# Plot ROC curve for each class
plt.figure(figsize=(10, 8))
for i in range(n_classes):
    plt.plot(fpr[i], tpr[i], label=f'Class {i} (AUC = {roc_auc[i]:.2f})')

plt.plot([0, 1], [0, 1], 'k--') # Plot diagonal line
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Each Class')
plt.legend(loc="lower right")
plt.show()
```



RESULTS

After preprocessing the dataset and training the machine learning model, we obtained promising results in predicting hotel listing discounts. The following key findings summarize the results of our analysis:

- **Regression Model Performance:** The developed regression model demonstrated robust performance in predicting the percentage of discount offered by hotel listings. Evaluation metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) indicated that the model's predictions were close to the actual discount values, suggesting high predictive accuracy.
- **Feature Importance:** Analysis of feature importance revealed that certain variables, such as price and rating, had a significant impact on determining the level of discount offered by hotels. Location also played a crucial role, with discounts varying across different cities or regions.
- **Generalization Ability:** The trained model exhibited good generalization ability, as it was able to accurately predict discounts for new hotel listings not present in the training data. This indicates that the model can provide reliable predictions for a wide range of accommodation options.
- **Insights for Travelers:** The predictive model can serve as a valuable tool for travelers seeking cost-effective accommodation options. By analyzing the predicted discounts for different hotels, travelers can make informed decisions and potentially save on accommodation costs during their trips.
- **Implications for Hotel Owners:** Hotel owners can leverage the insights provided by the model to optimize their pricing strategies and attract more customers. By understanding the factors influencing discount levels, hoteliers can adjust their pricing policies to maximize occupancy rates and revenue.

CONCLUSION

In conclusion, this research demonstrates the effectiveness of machine learning in predicting hotel listing discounts based on various factors such as location, rating, and price. By leveraging a dataset scraped from the official website of Trivago and employing regression modeling techniques, we were able to develop a predictive model with high accuracy in estimating discount percentages for hotel listings.

The findings of this study have several implications for both travelers and hotel owners. For travelers, the predictive model serves as a valuable tool for identifying cost-effective accommodation options and making informed booking decisions. By considering predicted discounts along with other factors such as location and amenities, travelers can optimize their accommodation choices and potentially save on lodging expenses during their trips.

For hotel owners, the insights provided by the predictive model offer valuable guidance in pricing optimization and revenue management. By understanding the impact of factors such as price competitiveness, rating, and geographical location on discount levels, hoteliers can adjust their pricing strategies to attract more customers and improve overall revenue generation.

Furthermore, this research contributes to the growing body of literature on data-driven approaches in the hospitality industry. By demonstrating the feasibility and effectiveness of machine learning in predicting hotel discounts, this study underscores the importance of leveraging data analytics to enhance decision-making processes and improve business outcomes in the hospitality sector.

REFERENCES

- “How to use Selenium to web-scrape with example”, Bryan Pfalzgraf, Published in Towards Data Science, April 3 2020
Link: <https://towardsdatascience.com/how-to-use-selenium-to-web-scrape-with-example-80f9b23a843a>
- “Understanding Logistic Regression in Python Tutorial”, Updated December 2019, datacamp.com
Link: <https://www.datacamp.com/tutorial/understanding-logistic-regression-python>
- “Guide to AUC-ROC curve in Machine Learning”, Aniruddha Bhandari, 23 April 2024, Analytics Vidhya
Link: <https://www.analyticsvidhya.com/blog/2020/06/auc-roc-curve-machine-learning/>