

# **Capstone Project**

## **EDA on Airbnb Booking Analysis**

**By : G MOHAMED LUQMAN**



WHAT WE ARE  
TALKING ABOUT ?

## INTRODUCTION

Airbnb is an online platform that connects hosts renting out space in their homes with guests seeking lodging for generally cheaper prices than a hotel.

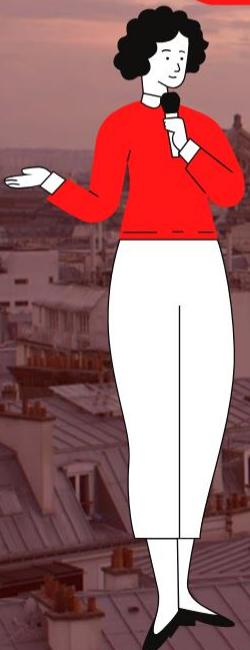
It has millions of listing, which generate alots of data. We are analyzing these data for making business decision, for looking best room type, etc.,



# Agenda

- 
- 1 Relationship analysis between variable
  - 2 Finding Top Neighborhood Group
  - 3 Finding total count of each room types
  - 4 Finding Most popular Neighborhood
  - 4 Types of room by neighbourhood regions
  - 6 Average price for each room type
  - 7 Rooms availability in different areas
  - 8 Price Analysis
  - 9 Listing Analysis
  - 10 Availability Analysis
  - 11 Profitability Analysis
  - 12 Visual Map Analysis and WordCount Analysis

# CONTENT



- 1 Data Exploration & variable Identification
- 2 Descriptive Statistical Analysis
- 3 Data Wrangling / Data Cleaning
- 4 Detecting & Handling Outliers
- 5 Data Visualization
- 6 Exploratory Data Analysis
- 7 Conclusion

# Understanding about the Dataset

<b>1. id</b>	Unique Listing Id	<b>9. Room_Type</b>	Listing space types
<b>2. Name</b>	Name of the Property	<b>10. Price</b>	Price in Dollars
<b>3. Host_id</b>	Unique Id for each listed host	<b>11. Minimum_nights</b>	Minimum nights required to stay
<b>4. Host_name</b>	Name of the host	<b>12. Number_of_reviews</b>	No. of reviews written for the listing
<b>5. Neighbourhood_Group</b>	Location	<b>13. Last_review</b>	Last reviewed date for the listing
<b>6. Neighbourhood</b>	Area	<b>14. Reviews_Per_Month</b>	Total review per month for the listing
<b>7. Latitude</b>	Latitude Coordinates	<b>15. Calculated_host_listings_count</b>	Total no of listing against the host id
<b>8. Longitude</b>	Longitude Coordinates	<b>16. Availability_365</b>	Number of days when listing is available for booking

# General Information about the Data

**Size of the dataset**

**782320**

**Shape of the Dataset**

**48895 (Rows)  
&  
16 (Columns)**

**Numerical Column  
(Quantitative Feature)**

**10**

**Categorical Column  
(Quanlitative Feature)**

**6**

# Descriptive Summary Statistics

	<b>id</b>	<b>host_id</b>	<b>latitude</b>	<b>longitude</b>	<b>price</b>	<b>minimum_nights</b>	<b>number_of_reviews</b>	<b>reviews_per_month</b>	<b>calculated_host_listings_count</b>	<b>availability_365</b>
<b>count</b>	4.889500e+04	4.889500e+04	48895.000000	48895.000000	48895.000000	48895.000000	48895.000000	38843.000000	48895.000000	48895.000000
<b>mean</b>	1.901714e+07	6.762001e+07	40.728949	-73.952170	152.720687	7.029962	23.274466	1.373221	7.143982	112.781327
<b>std</b>	1.098311e+07	7.861097e+07	0.054530	0.046157	240.154170	20.510550	44.550582	1.680442	32.952519	131.622289
<b>min</b>	2.539000e+03	2.438000e+03	40.499790	-74.244420	0.000000	1.000000	0.000000	0.010000	1.000000	0.000000
<b>25%</b>	9.471945e+06	7.822033e+06	40.690100	-73.983070	69.000000	1.000000	1.000000	0.190000	1.000000	0.000000
<b>50%</b>	1.967728e+07	3.079382e+07	40.723070	-73.955680	106.000000	3.000000	5.000000	0.720000	1.000000	45.000000
<b>75%</b>	2.915218e+07	1.074344e+08	40.763115	-73.936275	175.000000	5.000000	24.000000	2.020000	2.000000	227.000000
<b>max</b>	3.648724e+07	2.743213e+08	40.913060	-73.712990	10000.000000	1250.000000	629.000000	58.500000	327.000000	365.000000

## Observations:-

- **reviews\_per\_month** column has few missing values.
- Minimum value in **price** column is 0\$ and max is 10,000\$. (0\$ means free airbnbs, impossible, so bad data)
- There is/are properties whose **availability\_365** for booking is 0 day. (explored later)
- Maximum value for **minimum\_nights** column is 1250. (explored later)

```
[ ] # A function for price_correction
def price_correction(x):
    if x==0:
        return 100
    else:
        return x

[ ] df['price']=df['price'].apply(price_correction)

[ ] df['price'].apply(price_correction)

0      149
1      225
2      150
3       89
4       80
...
48890     70
48891     40
48892    115
48893     55
48894     90
Name: price, Length: 48895, dtype: int64
```

Minimum price is zero  
which is not possible. we  
should replace price value  
as 10 or 100



	<b>id</b>	<b>host_id</b>	<b>latitude</b>	<b>longitude</b>	<b>price</b>	<b>minimum_nights</b>	<b>number_of_reviews</b>	<b>reviews_per_month</b>	<b>calculated_host_listings_count</b>	<b>availability_365</b>
<b>count</b>	4 889500e+04	4 889500e+04	48895.000000	48895.000000	48895.000000	48895.000000	48895.000000	38843.000000	48895.000000	48895.000000
<b>mean</b>	1.901714e+07	6.762001e+07	40.728949	-73.952170	152.743184	7.029962	23.274466	1.373221	7.143982	112.781327
<b>std</b>	1.098311e+07	7.861097e+07	0.054530	0.046157	240.144546	20.510550	44.550582	1.680442	32.952519	131.622289
<b>min</b>	2.539000e+03	2.438000e+03	40.499790	-74.244420	10.000000	1.000000	0.000000	0.010000	1.000000	0.000000
<b>25%</b>	9.471945e+06	7.822033e+06	40.690100	-73.983070	69.000000	1.000000	1.000000	0.190000	1.000000	0.000000
<b>50%</b>	1.967728e+07	3.079382e+07	40.723070	-73.955680	106.000000	3.000000	5.000000	0.720000	1.000000	45.000000
<b>75%</b>	2.915218e+07	1.074344e+08	40.763115	-73.936275	175.000000	5.000000	24.000000	2.020000	2.000000	227.000000
<b>max</b>	3.648724e+07	2.743213e+08	40.913060	-73.712990	10000.000000	1250.000000	629.000000	58.500000	327.000000	365.000000

	<b>id</b>	<b>host_id</b>	<b>latitude</b>	<b>longitude</b>	<b>price</b>	<b>minimum_nights</b>	<b>number_of_reviews</b>	<b>reviews_per_month</b>	<b>calculated_host_listings_count</b>	<b>availability_365</b>
<b>count</b>	4.889500e+04	4.889500e+04	48895.000000	48895.000000	48895.000000	48895.000000	48895.000000	38843.000000	48895.000000	48895.000000
<b>mean</b>	1.901714e+07	6.762001e+07	40.728949	-73.952170	152.743184	6.942980	23.274466	1.373221	7.143982	112.781327
<b>std</b>	1.098311e+07	7.861097e+07	0.054530	0.046157	240.144546	17.530294	44.550582	1.680442	32.952519	131.622289
<b>min</b>	2.539000e+03	2.438000e+03	40.499790	-74.244420	10.000000	1.000000	0.000000	0.010000	1.000000	0.000000
<b>25%</b>	9.471945e+06	7.822033e+06	40.690100	-73.983070	69.000000	1.000000	1.000000	0.190000	1.000000	0.000000
<b>50%</b>	1.967728e+07	3.079382e+07	40.723070	-73.955680	106.000000	3.000000	5.000000	0.720000	1.000000	45.000000
<b>75%</b>	2.915218e+07	1.074344e+08	40.763115	-73.936275	175.000000	5.000000	24.000000	2.020000	2.000000	227.000000
<b>max</b>	3.648724e+07	2.743213e+08	40.913060	-73.712990	10000.000000	365.000000	629.000000	58.500000	327.000000	365.000000

maximum night should be  
365 days not as 1250  
days. we should fix it.



```
def minimum_nights_correction(x):
    if x>365:
        return 365
    else:
        return x

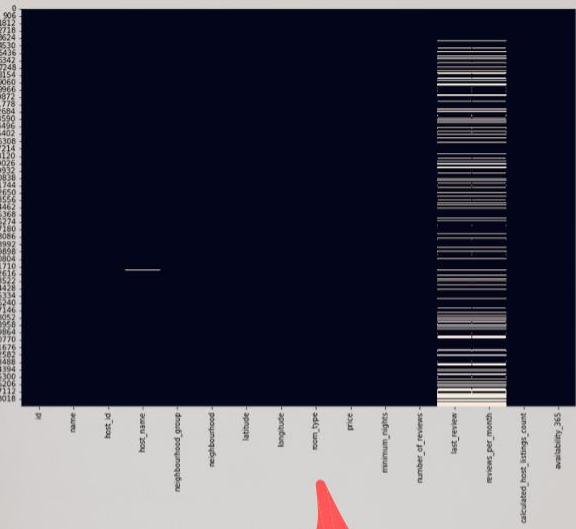
df['minimum_nights']=df['minimum_nights'].apply(minimum_nights_correction)

df[df['minimum_nights']==365]
```

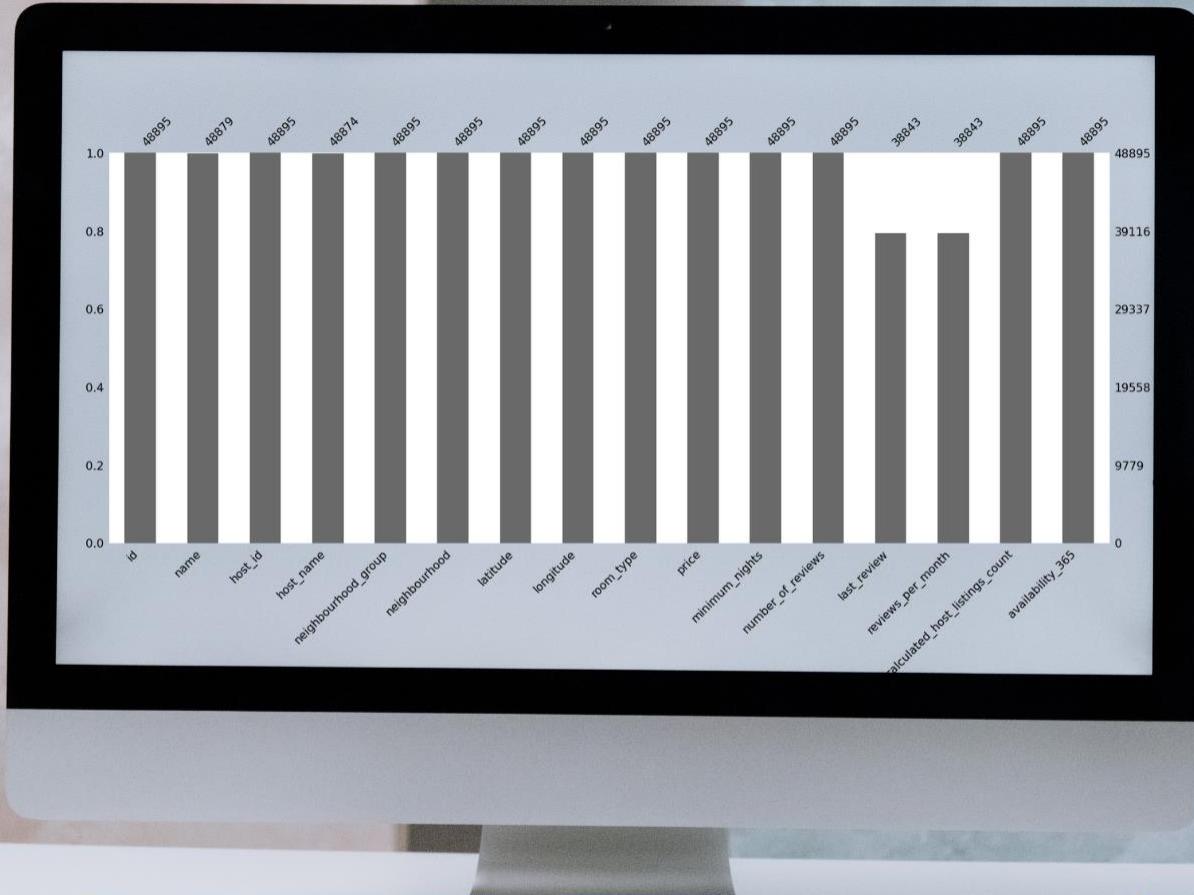
## NULL Values in the Dataset

```
[ ] # Checking for missing values
df.isnull().sum()

id          0
name         16
host_id      0
host_name    21
neighbourhood_group  0
neighbourhood   0
latitude        0
longitude       0
room_type       0
price          0
minimum_nights  0
number_of_reviews 0
last_review    10052
reviews_per_month 10052
calculated_host_listings_count 0
availability_365 0
dtype: int64
```



Visualization of Null values using Heatmap



using missingno plot can display  
count of non - null values

# Handling Missing Value

```
def impute_median(series):
    return series.fillna(series.median())

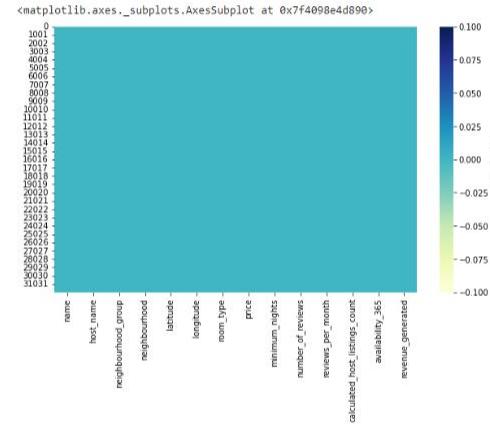
df.reviews_per_month=df["reviews_per_month"].transform(impute_median)

df['host_name'].fillna("unavailable",inplace = True)
df['name'].fillna("unavailable",inplace = True)

df.isnull().sum()

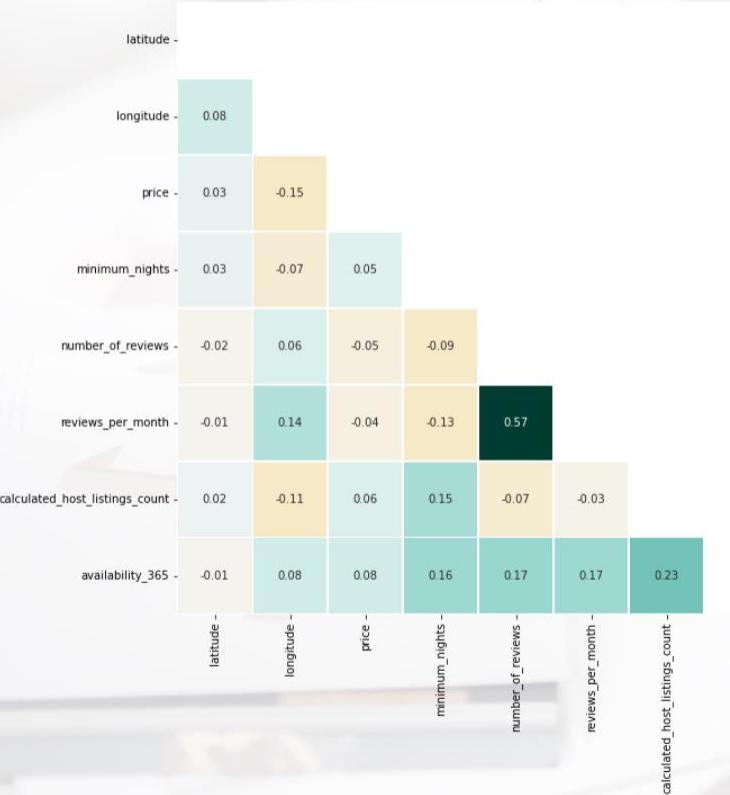
name          0
host_name     0
neighbourhood_group 0
neighbourhood 0
latitude       0
longitude      0
room_type      0
price          0
minimum_nights 0
number_of_reviews 0
reviews_per_month 0
calculated_host_listings_count 0
availability_365 0
dtype: int64

plt.figure(figsize=(10,6))
sns.heatmap(df.isna(),cmap="YlGnBu",cbar_kws={'label': 'Missing Data'})
```



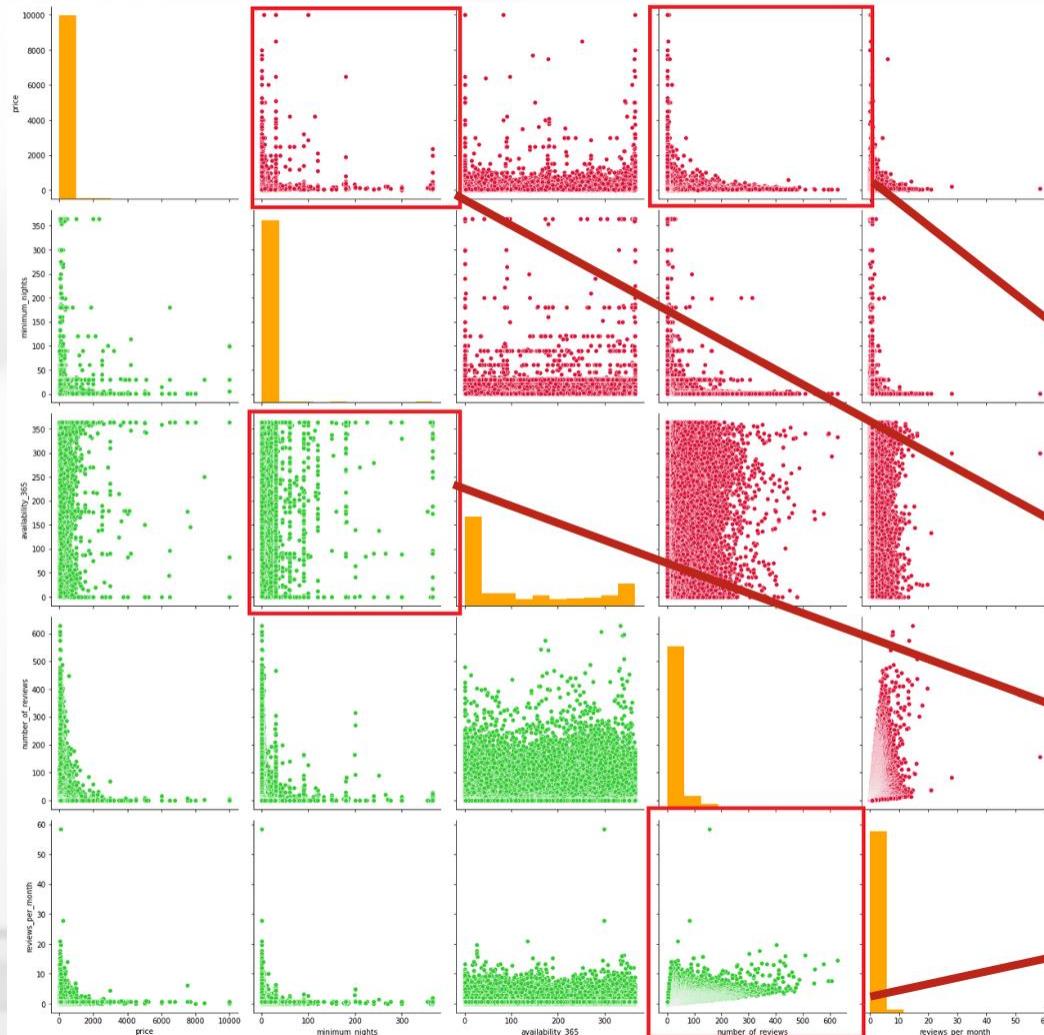
# Correlation

Feature-correlation (pearson)



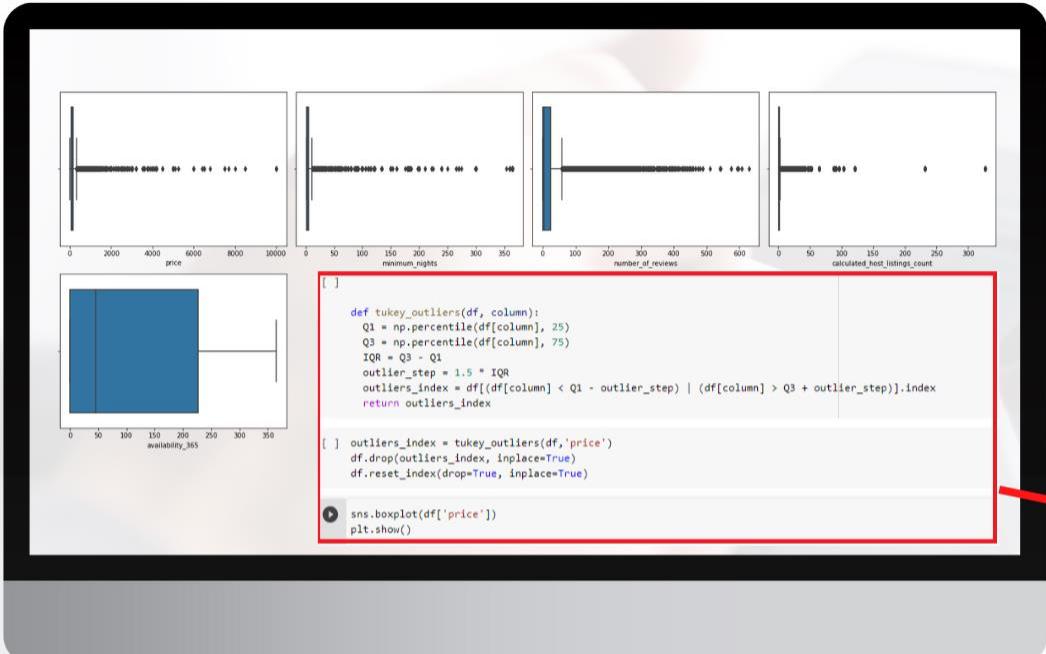
Categorical data plot



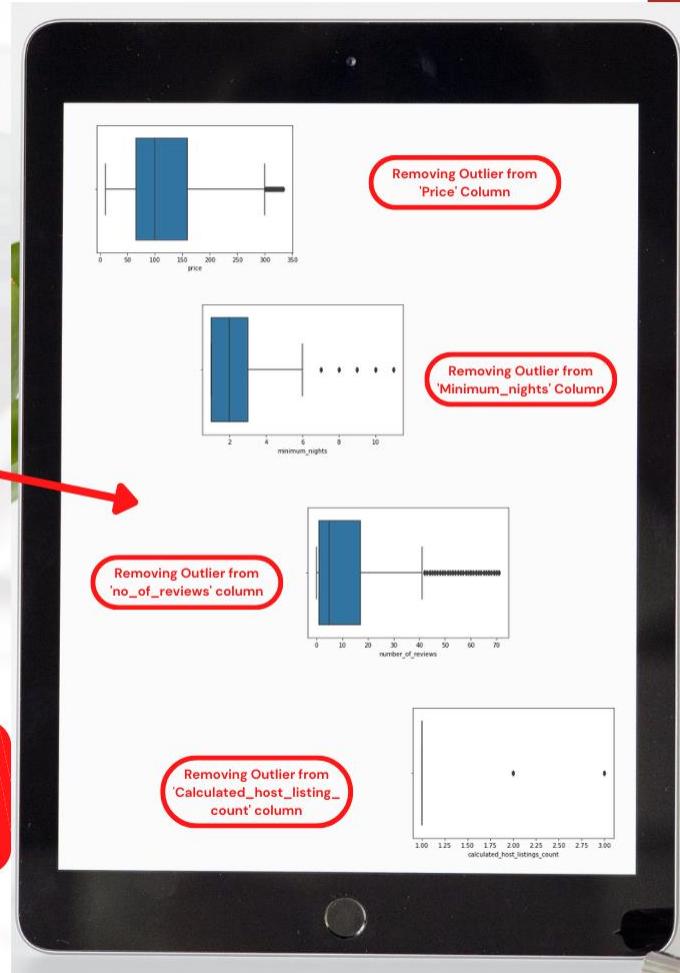


## Relationship analysis between variable (pairplot)

- 1** There are more number of reviews where price is low
- 2** As minimum nights of booking increase price decreases significantly.
- 3** Unexpected minimum nights at 0 availability
- 4** Number of reviews and review per month follows linear relationship between them. Some plot has a extreme points



```
[ ]  
def tukey_outliers(df, column):  
    Q1 = np.percentile(df[column], 25)  
    Q3 = np.percentile(df[column], 75)  
    IQR = Q3 - Q1  
    outlier_step = 1.5 * IQR  
    outliers_index = df[(df[column] < Q1 - outlier_step) | (df[column] > Q3 + outlier_step)].index  
    return outliers_index  
  
[ ] outliers_index = tukey_outliers(df,'price')  
df.drop(outliers_index, inplace=True)  
df.reset_index(drop=True, inplace=True)  
  
sns.boxplot(df['price'])  
plt.show()
```



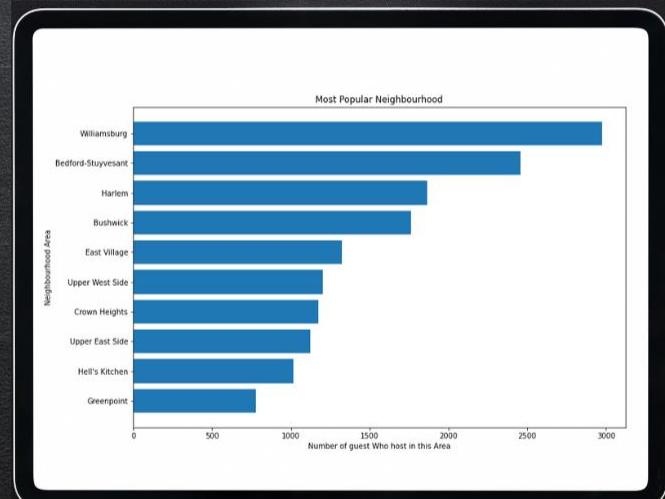
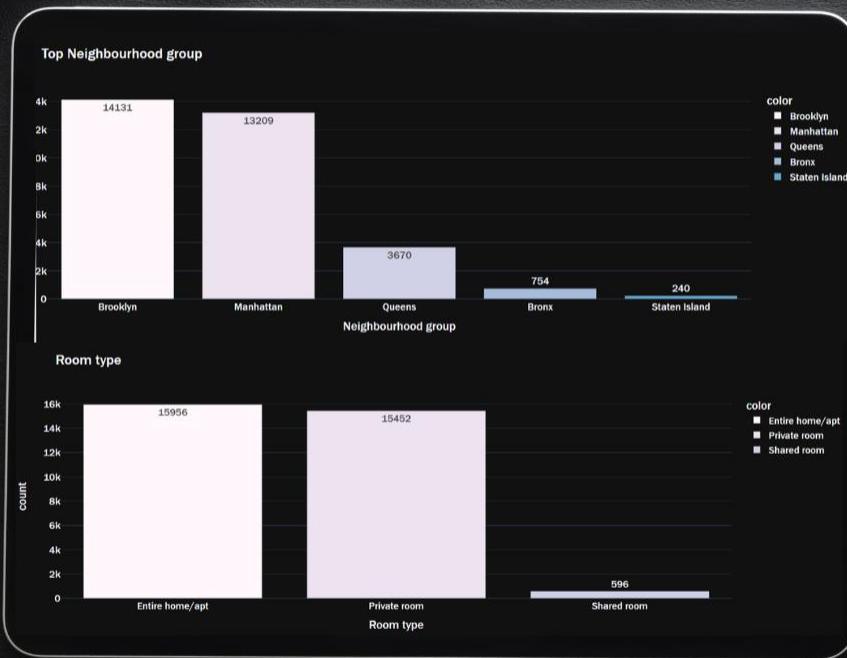
Removing  
Outliers



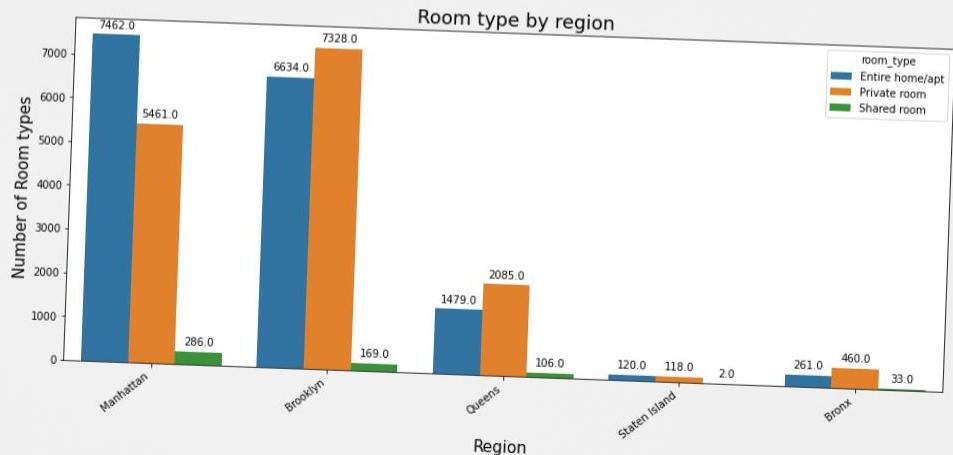
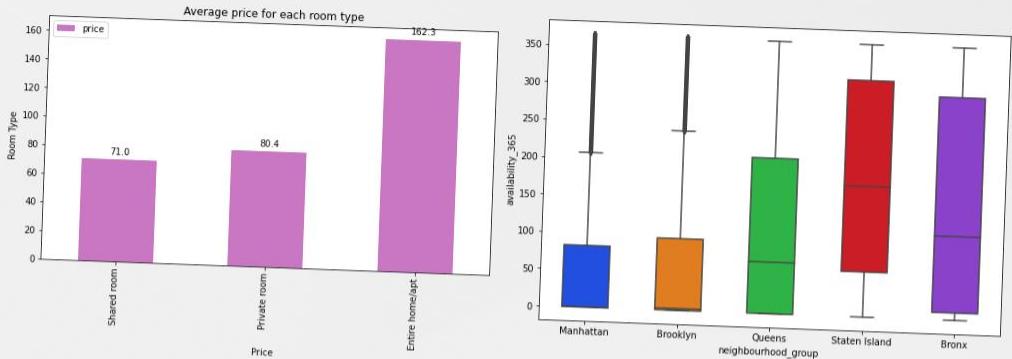
DATA  
VISUALIZATION



# Univariate Analysis



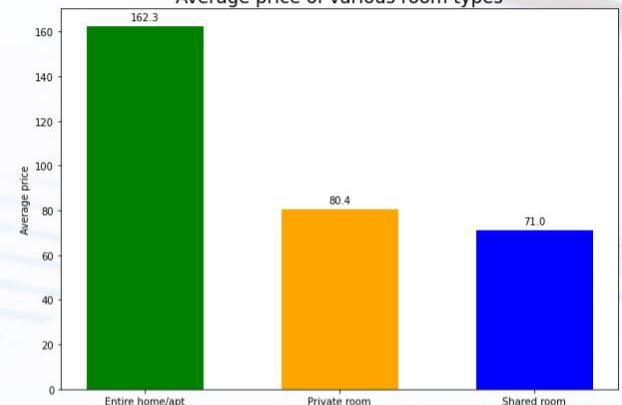
## Bivariate Analysis



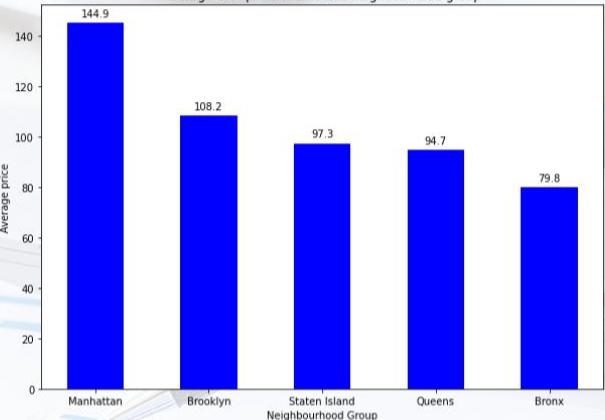
# EXPLORATORY DATA ANALYSIS

## Price Analysis

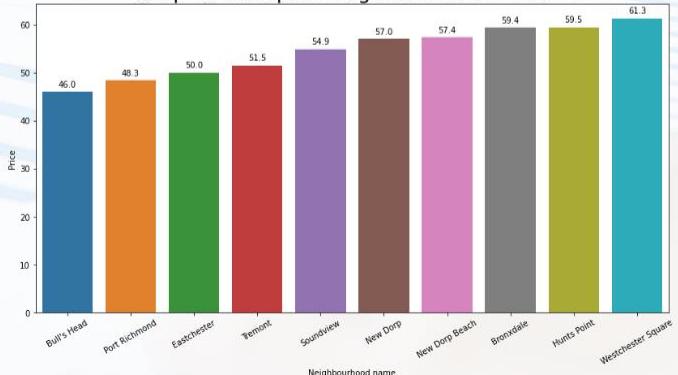
Average price of various room types



Average cost price of different neighbourhood group



Top 10 Cheapest neighbourhood in NYC



## Observation:

Order of costliest neighbourhood group on basis of average price and price distribution => Manhattan > Brooklyn > Staten Island > Queens > Bronx.

Staten Island was having least no. of listing but it is not the cheapest.

For buying a property to get business with airbnb, Brooklyen and queens are most preferred neighbourhood group because of less saturation level compaired to Manhattan

## EXPLORATORY DATA ANALYSIS

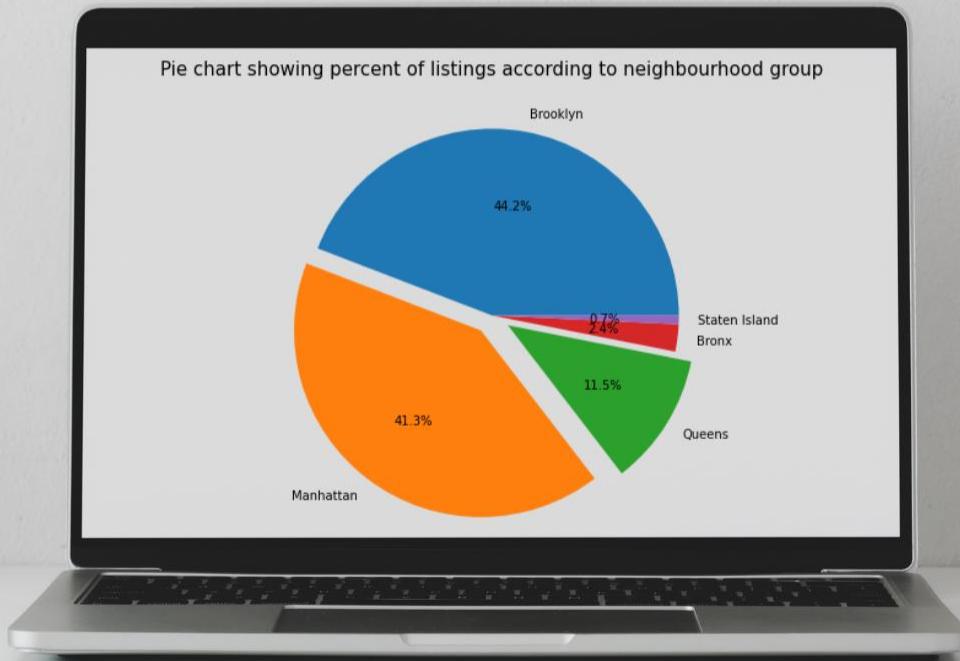
### Result:

1. Manhatten & Brooklyn are having high no. of listing.
2. Staten island and Bronx have low no. of listing.

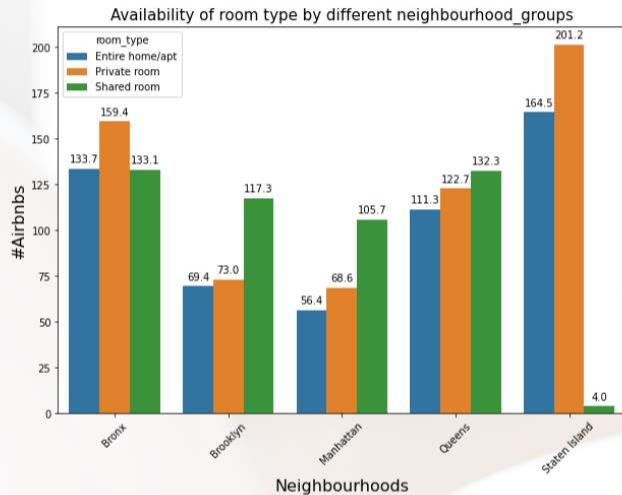
### Inference:

If we want to do advertisement or marketing for selling flats we should focus on Manhattan and Brooklyn.

### Listing Analysis



## EXPLORATORY DATA ANALYSIS



## Availability analysis

### Result:

1. Private room has highest mean availability..
2. Entire home has least mean availability..

### Inference:

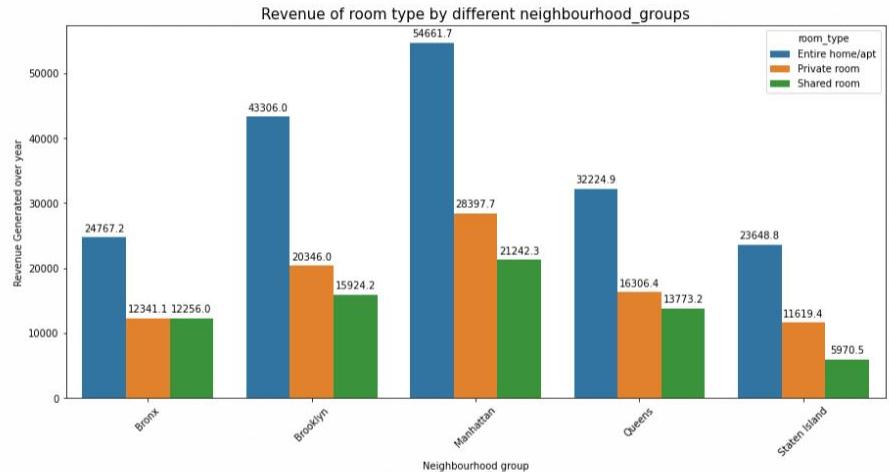
Manhattan and Brooklyn have less availability than compared to other neighbourhood groups which is good for the host having these properties.

# Profitability Analysis

## EXPLORATORY DATA ANALYSIS

### Observation:

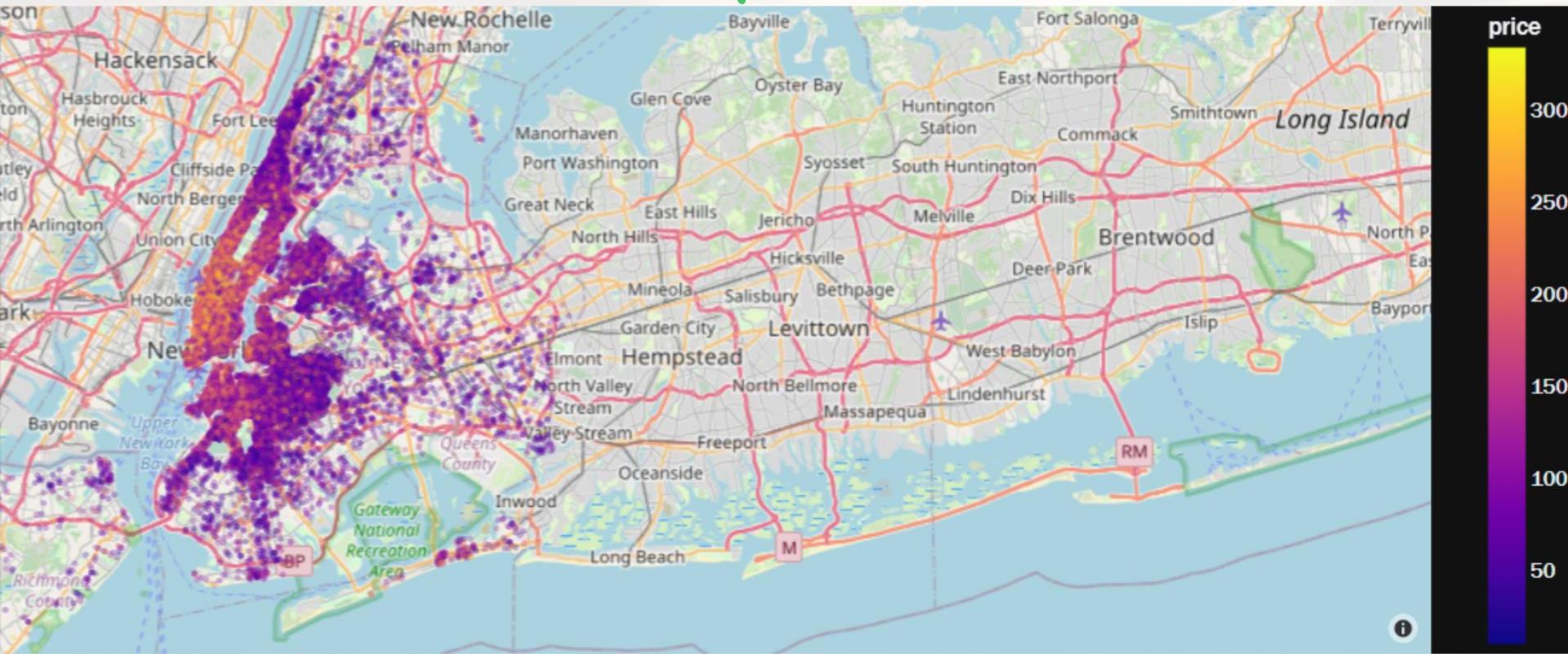
As we can see that in every neighbourhood "Entire home" has generated most revenue. So buying a property "Entire home" and renting it, is a profitable business irrespective of neighbourhood type.



## Visualization on Maps using Latitude and longitude



# Visual Analysis using Satellite Map



## Maps of Room Type



# Word Cloud of the Most Keyword Used



# CONCLUSION

- 
- 1 Entire home/apt is highly expensive.
- 2 Manhattan living cost is highest, Bronx living cost is lowest
- 3 Cheapest neighbourhood is Bulls head.
- 4 Manhattan have the highest no. of listing.
- 5 Private room has the highest availability;
- 6 the Entire home has least availability
- 7 Manhattan, Brookyln & some parts of Queens have a high traffic of Booking
- 8 Most visitors don't prefer shared rooms, They tends to visit private or entire home/ apt.
- 9 Williamsburg is the neighbourhood with highest number of listings.
- 10 Private rooms,Apartment,Brooklyn,Manhattan are most number time used keywords
- 11 Queens is the neighbourhood which has received maximum number of reviews, bronx was least.