

# **Airbnb NYC- Data-Driven Analysis**

## **Methodology**

By:

Sagar Barge, Vunna Praveen Kumar, Satbir Kaur

# Methodology

## 1. Understanding problem & objective:

Airbnb has experienced a significant decline in revenue over recent months, likely due to reduced travel demand. So, Airbnb wants to find solutions for this with help of listing data.

## 2. Understanding the Data-

This data includes all Airbnb listings in New York City. Contains details such as host information, property type, price, location, availability, and customer reviews.

Data Dictionary

Column	Description
id	listing ID
name	name of the listing
host_id	host ID
host_name	name of the host
neighbourhood_group	location
neighbourhood	area
latitude	latitude coordinates
longitude	longitude coordinates
room_type	listing space type
price	
minimum_nights	amount of nights minimum
number_of_reviews	number of reviews
last_review	latest review
reviews_per_month	number of reviews per month
calculated_host_listings_count	amount of listing per host
availability_365	number of days when listing is available for booking

## 3. Discussing about Tools required:

- After discussion with team members we decided to different software to achieve this objective.
- Tools: Python, Excel, PowerBI, Tableau, QGIS.

## 4. Import required libraries-

### Import Liabraries

```
In [1]: 1 # Import Liabraries
        2 import numpy as np
        3 import pandas as pd
        4 import matplotlib.pyplot as plt
        5 import seaborn as sns
        6 import geopandas as gpd
        7 from shapely.geometry import Point
        8 import warnings
        9 warnings.filterwarnings ('ignore')
```

## 5. Checking data type, and basic statistics to get an overview of data:

- `.shape, df.info(), df.describe()`

## 6. Checking duplicate values:

```
1 # Find duplicates
2 df.duplicated().sum()
0
```

## 7. Checking null values :

```
1 # Checking null values
2
3 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
```

## 8. Filled Null values :

- Filled 'name', 'host\_name' as 'Unknown'. This column is useful for finding top listing names.
- Filled Null values as 0 >> Because “number\_of\_reviews” column having 0 value, so “reviews\_per\_month” will be 0.
- Changing data format to date

```
1 # Filled Null values as Unknown
2 df[['name', 'host_name']] = df[['name', 'host_name']].fillna('Unknown')

1 # Filled Null values as 0 >> Because number_of_reviews column having 0 value, so reviews_per_month will be 0.
2 df['reviews_per_month'] = df['reviews_per_month'].fillna(0)

1 # Changing data format to date
2 df['last_review'] = pd.to_datetime(df['last_review'])
3
```

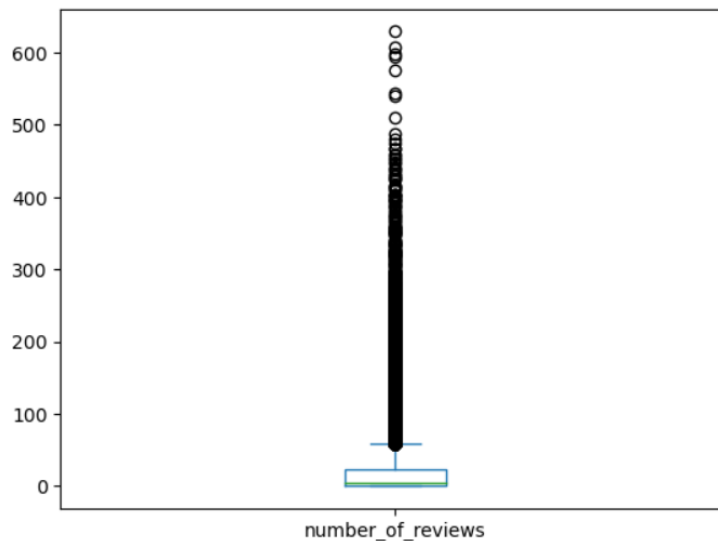
## 9. Checking outlier:

- For checking outlier used box plot & value count

Example-

```
1 df.number_of_reviews.plot.box()
```

&lt;Axes: &gt;



- “minimum nights” capped at 500, because after 500, there were only a few values.

```
1 # Capping at 500
2
3 df['minimum_nights'] = df['minimum_nights'].where(df['minimum_nights'] < 500, 500)
4 df.minimum_nights.plot.box()
```

## 10. Dropping Unnecessary columns

```
1 # Dropping Unnecessary columns
2 df = df.drop(['last review', 'id'], axis=1)
```

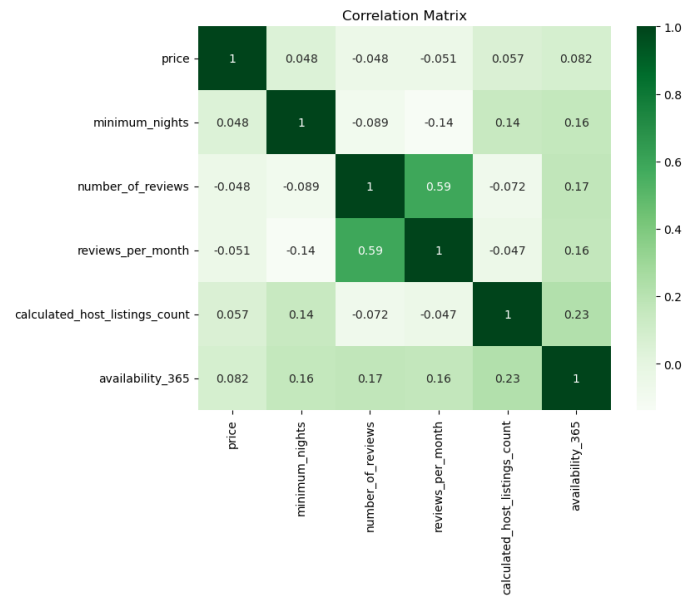
## 11. Creating new Column with Low review (0), to understand unpopular listings

```
1 # Adding Column with Low review (0)
2 df['Low_review'] = df['number_of_reviews']==0
3 df['Low_review 01'] = (df['number_of_reviews']==0).astype(int)
```

## 12. Export final dataframe to .csv for more analysis work >>>>>>>>>>

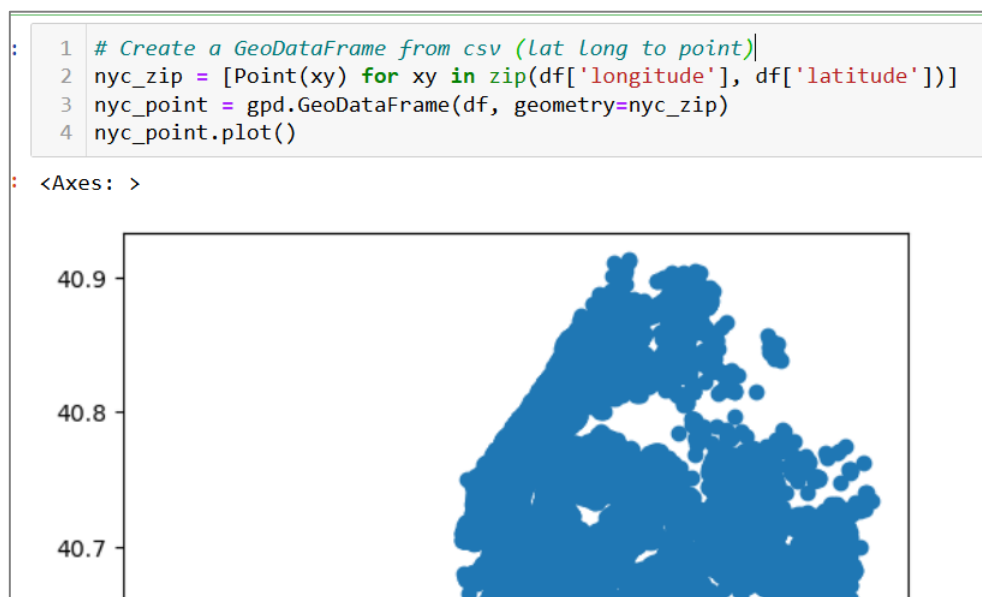
```
1 # Export dataframe to .csv
2
3 df.to_csv('Airbnb NYC final.csv')
```

### 13. Correlation Matrix to understand relation between different variables

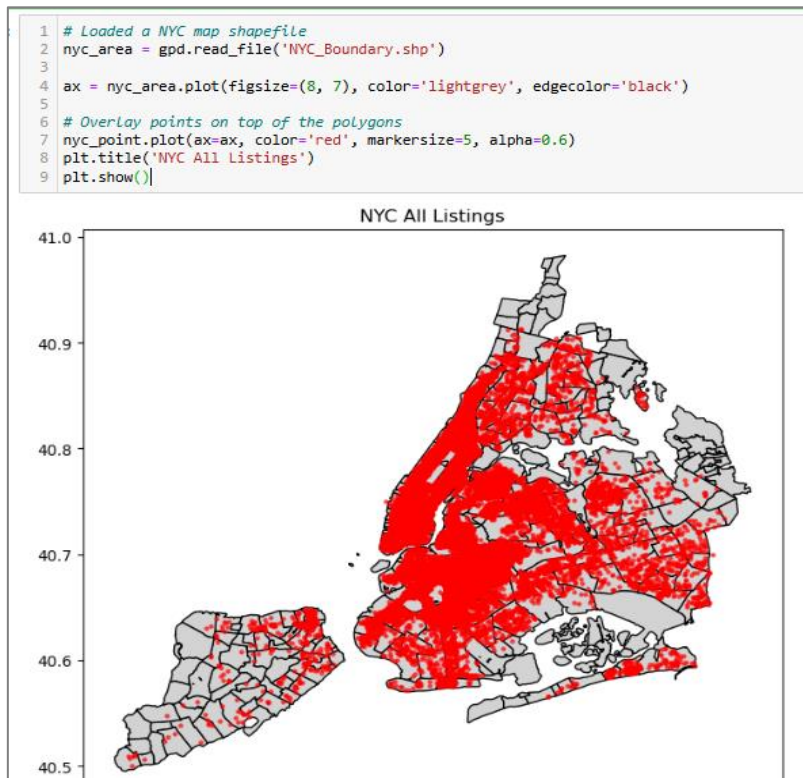


### 14. Hotspot map created using geopandas

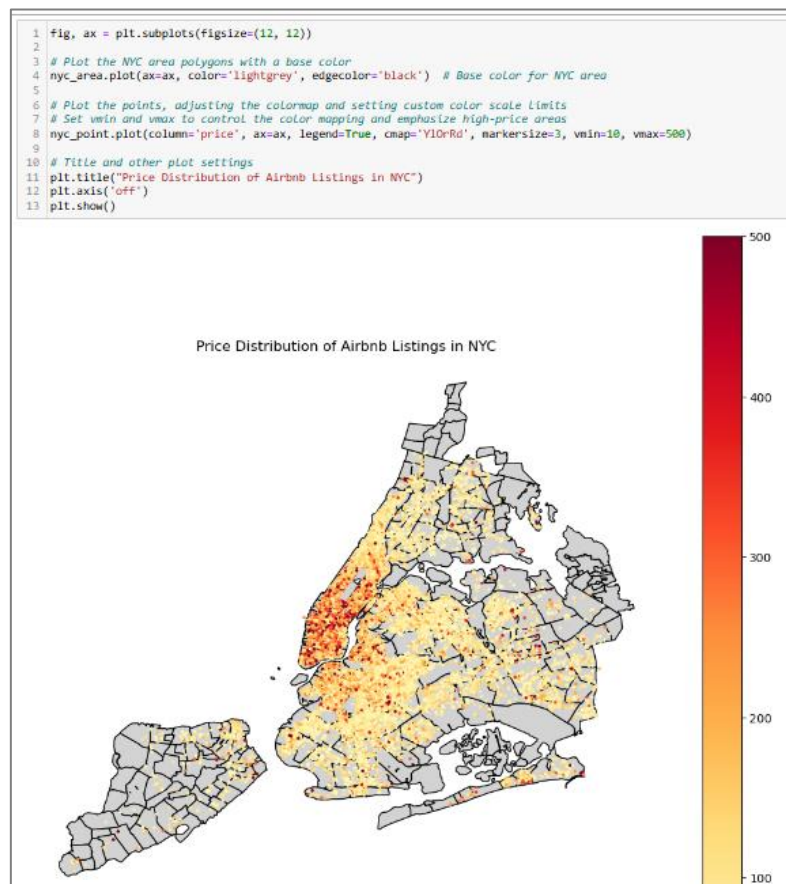
- Created a GeoDataFrame from .csv (lat long to point)



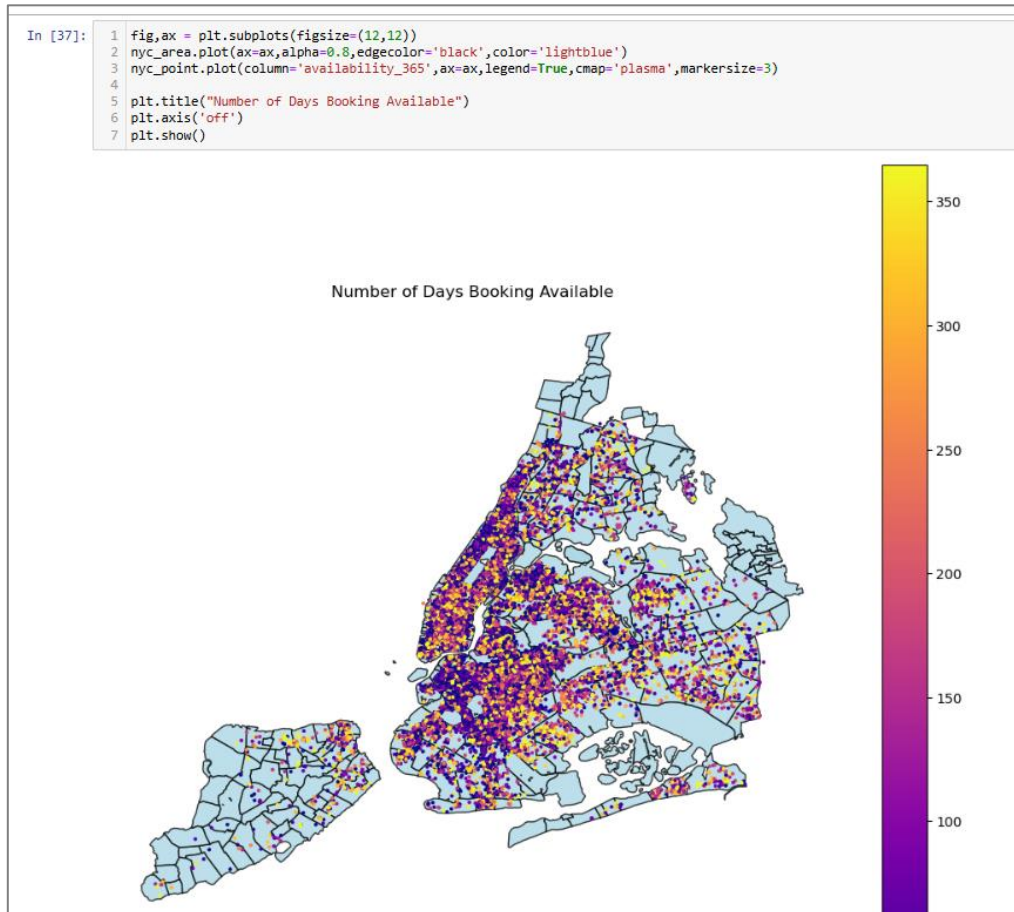
- Loaded NYC map shapefile over listing point



- Price Distribution

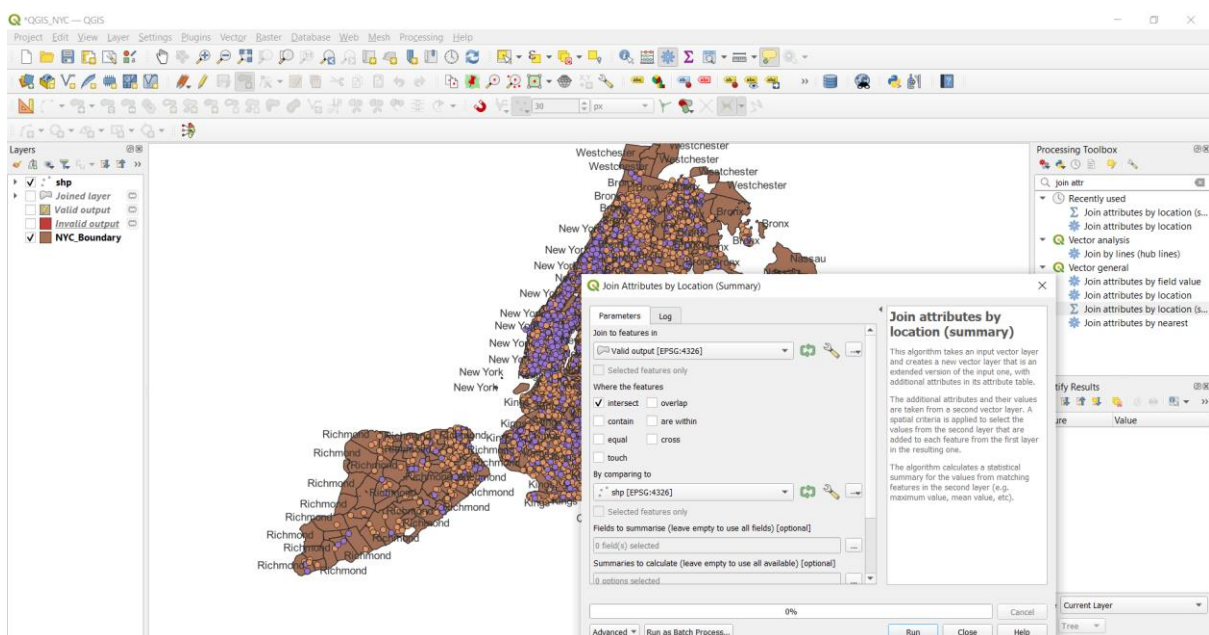


- Number of Days Booking Available



- Join attribute in QGIS (point and neighborhood boundary layer)

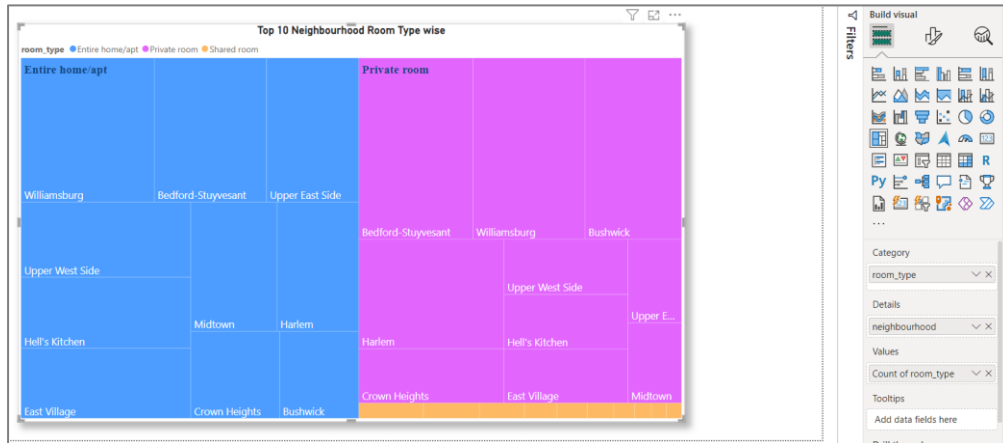
To identify Unpopular Properties, we created new columns as “Low\_review” where 1= 0 review & 0 = other reviews. we sorted out properties with No Reviews and showed them on a QGIS map.



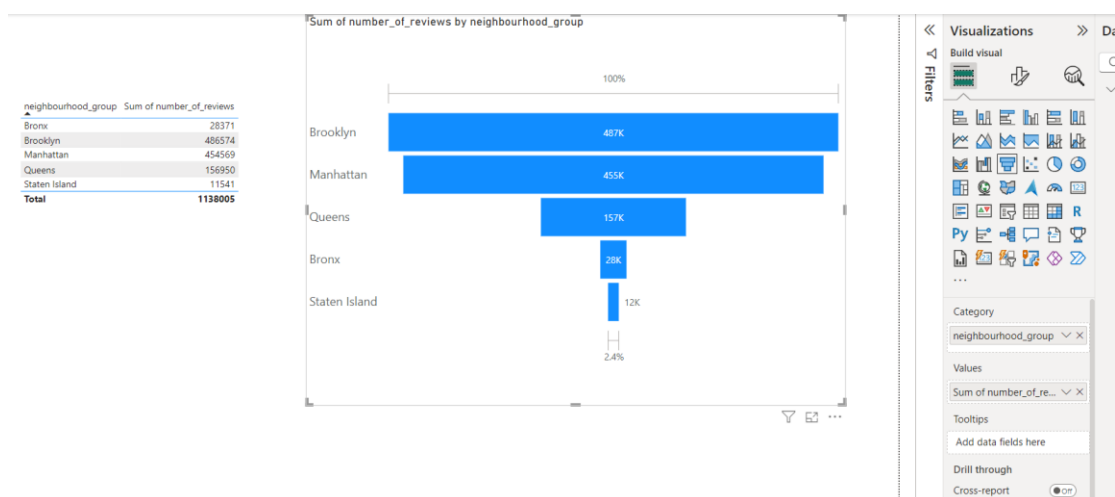
## Power BI, Tableau

### 15. Importing data in Power BI, Tableau & Analysis

- To analyze which type of host to acquire more and where, we analyzed the Top 10 neighborhoods with the type of rooms to count of listings. the following chart shows which listings are available and where they are.



- Neighborhoods group by count of listings, Here we understand listing count with their size of shape.

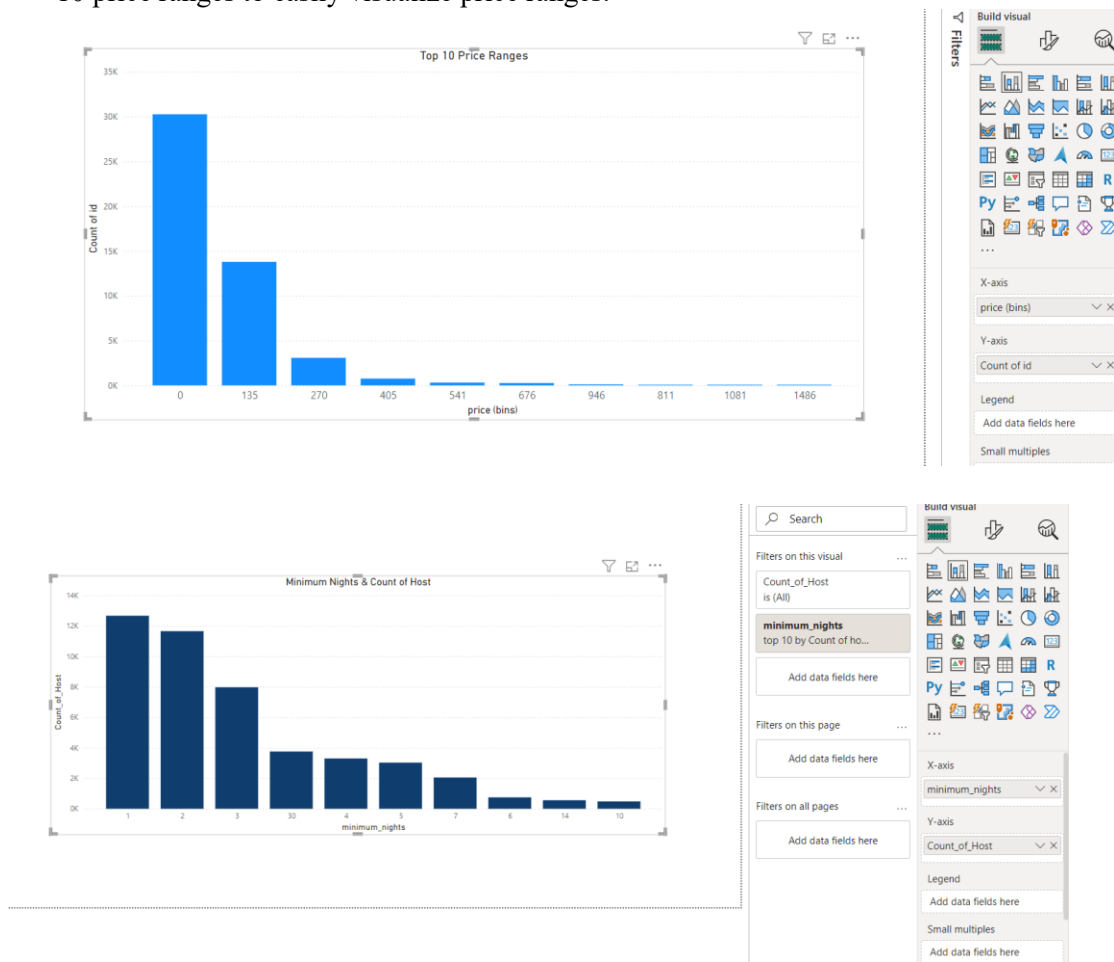


### Top 10 Neighborhoods by Sum of No. of reviews

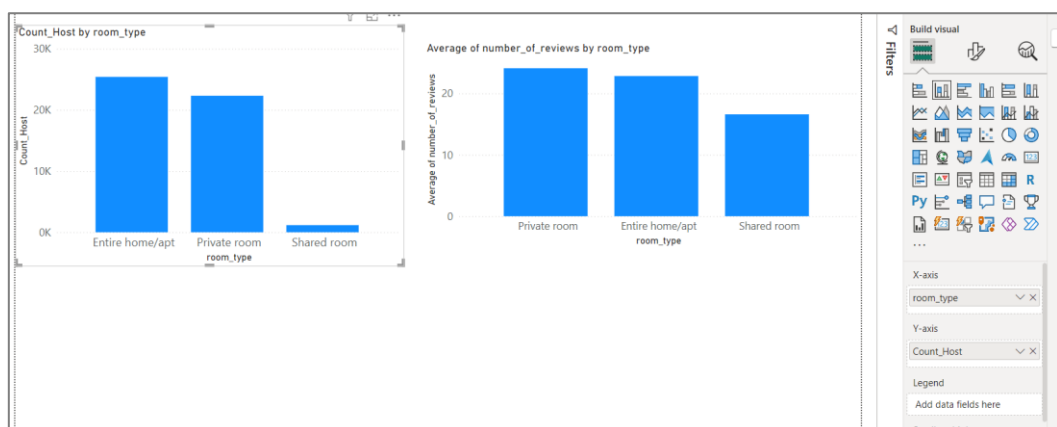




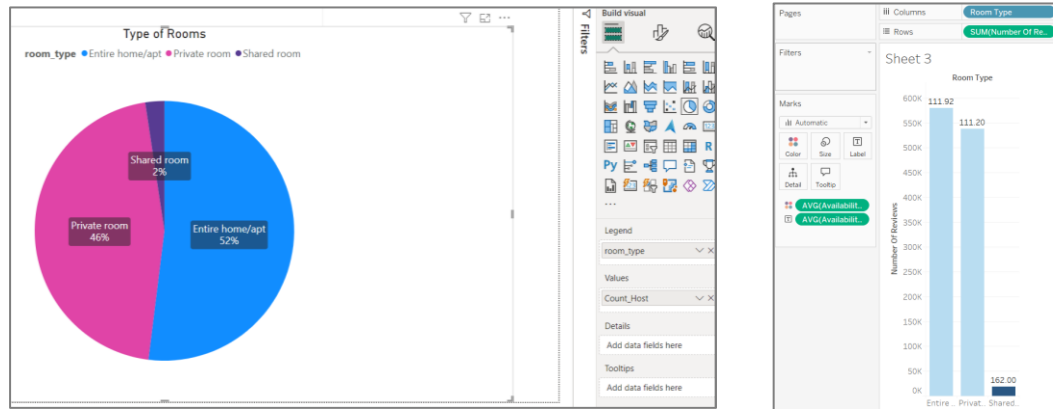
- To know the categorization of customers based on their preferences we created bins of price ranges & to know preferred stay for days. There was an outlier in the data so filtered as top 10 price ranges to easily visualize price ranges.



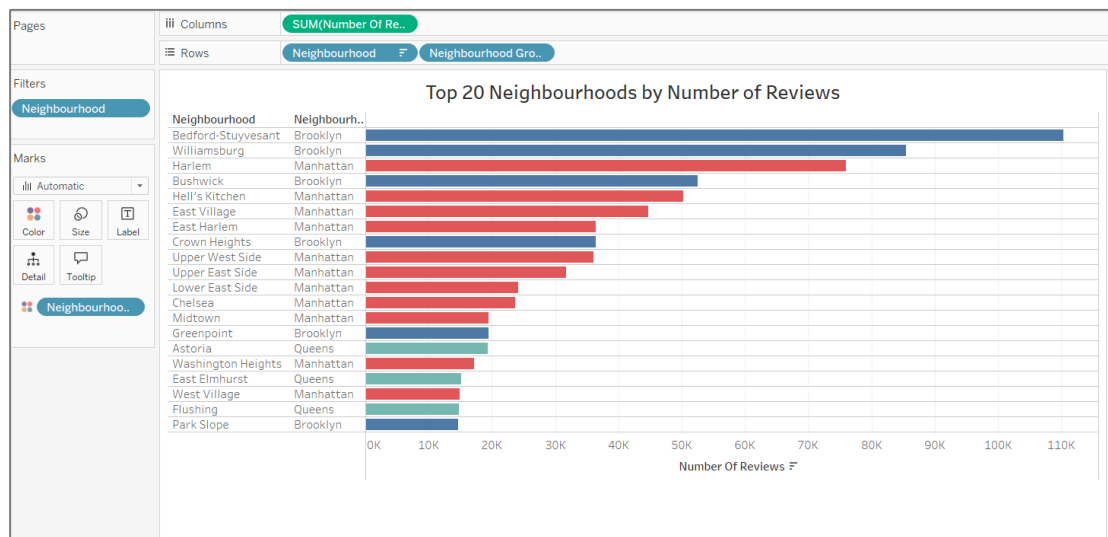
- To know The various kinds of properties that exist w.r.t. customer preferences, we calculated count of host by room type & Average review room type



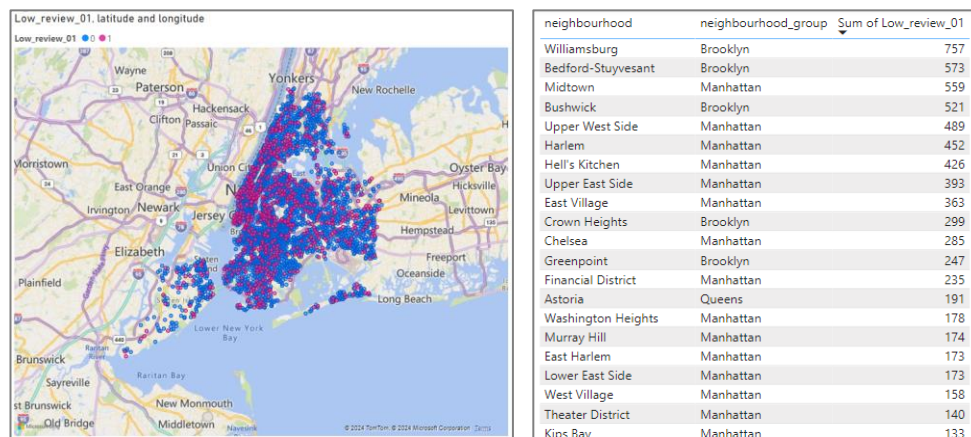
- To know adjustments in the existing properties to make it more customer-oriented in the first pie-chart share of different room type presented. In second plot bar chart created with room type in x axis, sum of number of reviews on y axis, the average availability given to text labels to bars & color shades according to average availability.



- When calculating the most popular localities in New York, we used a bar chart here, with the sum of the number of reviews in columns and the neighborhood and neighborhood group in rows. We then filtered the top 20 Properties by the number of reviews.



- To identify Unpopular Properties, we created new columns as “Low\_review” where 1= 0 review & 0 = other reviews. we sorted out properties with No Reviews and showed them on a location map. The map highlights unpopular properties. with this we presented Unpopular Neighbourhood Group, top 20 Unpopular Properties by



- Assumptions:**
  - Airbnb assumes that after covid-19 pandemic travel activity will increase.
  - Identified customer preferences using the number of reviews given by customers