# UNVEILING THE SECRETS OF AIRBNB IN NYC: DATA INSIGHTS

# AGENDA

Objective

Data life cycle

Analysis methods

Recommendations

Appendix:

- Data sources
- Data methodology
- Data model assumptions

# OBJECTIVE

To Conduct a thorough analysis of New York Airbnb Dataset.

Ask effective questions that can lead to data insights

Process, analyse and share findings by data visualization and statistical techniques

# DATA LIFE CYCLE

In the first phase the data captured and loaded into various environment.

Once data is cleaned, EDA is done and new features are created.

Then Meaningful insights are derived using various analytical methods.

# 1. Importing libraries and reading the data

```
In [1]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
```

```
In [2]:  inp0 = pd.read_csv(r'C:\Users\nyk\Downloads\AB_NYC_2019.csv')
         inp0.head(10)
```

Out[2]:

| | id | name | host_id | host_name | neighbourhood_group | neighbourhood | latitude | longitude | room_type | price | minimum_nights | number_of_revie |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2539 | Clean & quiet apt home by the park | 2787 | John | Brooklyn | Kensington | 40.64749 | -73.97237 | Private room | 149 | 1 | |
| 1 | 2595 | Skylit Midtown Castle | 2845 | Jennifer | Manhattan | Midtown | 40.75362 | -73.98377 | Entire home/apt | 225 | 1 | |
| 2 | 3647 | THE VILLAGE OF HARLEM....NEW YORK ! | 4632 | Elisabeth | Manhattan | Harlem | 40.80902 | -73.94190 | Private room | 150 | 3 | |
| 3 | 3831 | Cozy Entire Floor of Brownstone | 4869 | LisaRoxanne | Brooklyn | Clinton Hill | 40.68514 | -73.95976 | Entire home/apt | 89 | 1 | |
| 4 | 5022 | Entire Apt: Spacious Studio/Loft by central park | 7192 | Laura | Manhattan | East Harlem | 40.79851 | -73.94399 | Entire home/apt | 80 | 10 | |

# 2. Creating features

```python
def availability_365_categories_function(row):
    """
    Categorizes the "minimum_nights" column into 5 categories
    """
    if row <= 1:
        return 'very Low'
    elif row <= 100:
        return 'Low'
    elif row <= 200 :
        return 'Medium'
    elif (row <= 300):
        return 'High'
    else:
        return 'very High'
```

**2.2 categorizing the "minimum_nights" column into 5 categories**

```python
def minimum_night_categories_function(row):
    """
    Categorizes the "minimum_nights" column into 5 categories
    """
    if row <= 1:
        return 'very Low'
    elif row <= 3:
        return 'Low'
    elif row <= 5 :
        return 'Medium'
    elif (row <= 7):
        return 'High'
    else:
        return 'very High'
```

**2.3 categorizing the "number_of_reviews" column into 5 categories**

```python
def number_of_reviews_categories_function(row):
    """
    Categorizes the "number_of_reviews" column into 5 categories
    """
    if row <= 1:
        return 'very Low'
    elif row <= 5:
        return 'Low'
    elif row <= 10 :
        return 'Medium'
    elif (row <= 30):
        return 'High'
    else:
        return 'very High'
```

# 3. Fixing columns

```python
import warnings
warnings.filterwarnings("ignore")
inp0.last_review = pd.to_datetime(inp0.last_review)
inp0.last_review
```

```
0        2018-10-19
1        2019-05-21
2               NaT
3        2019-05-07
4        2018-11-19
            ...
48890           NaT
48891           NaT
48892           NaT
48893           NaT
48894           NaT
Name: last_review, Length: 48895, dtype: datetime64[ns]
```

```python
inp0.columns
```

```
Index(['id', 'name', 'host_id', 'host_name', 'neighbourhood_group',
       'neighbourhood', 'latitude', 'longitude', 'room_type', 'price',
       'minimum_nights', 'number_of_reviews', 'last_review',
       'reviews_per_month', 'calculated_host_listings_count',
       'availability_365', 'availability_365_categories',
       'minimum_night_categories', 'number_of_reviews_categories',
       'price_categories'],
      dtype='object')
```

# 4. Data types

```
: inp0.columns
```
```
: Index(['id', 'name', 'host_id', 'host_name', 'neighbourhood_group',
         'neighbourhood', 'latitude', 'longitude', 'room_type', 'price',
         'minimum_nights', 'number_of_reviews', 'last_review',
         'reviews_per_month', 'calculated_host_listings_count',
         'availability_365', 'availability_365_categories',
         'minimum_night_categories', 'number_of_reviews_categories',
         'price_categories'],
        dtype='object')
```

```
: # Categorical nominal
  categorical_columns = inp0.columns[[0,1,3,4,5,8,16,17,18,19]]
  categorical_columns
```
```
: Index(['id', 'name', 'host_name', 'neighbourhood_group', 'neighbourhood',
         'room_type', 'availability_365_categories', 'minimum_night_categories',
         'number_of_reviews_categories', 'price_categories'],
        dtype='object')
```

**4.2 Numerical**

```
numerical_columns = inp0.columns[[9,10,11,13,14,15]]
numerical_columns
```
```
Index(['price', 'minimum_nights', 'number_of_reviews', 'reviews_per_month',
       'calculated_host_listings_count', 'availability_365'],
      dtype='object')
```

```
inp0[numerical_columns].head()
```

|   | price | minimum_nights | number_of_reviews | reviews_per_month | calculated_host_listings_count | availability_365 |
|---|-------|----------------|-------------------|-------------------|--------------------------------|------------------|
| 0 | 149 | 1 | 9 | 0.21 | 6 | 365 |
| 1 | 225 | 1 | 45 | 0.38 | 2 | 355 |
| 2 | 150 | 3 | 0 | NaN | 1 | 365 |
| 3 | 89 | 1 | 270 | 4.64 | 1 | 194 |
| 4 | 80 | 10 | 9 | 0.10 | 1 | 0 |

**4.3 Coordinates and date**

```
coordinates = inp0.columns[[5,6,12]]
inp0[coordinates]
```

|   | neighbourhood | latitude | last_review |
|---|---------------|----------|-------------|
| 0 | Kensington | 40.64749 | 2018-10-19 |
| 1 | Midtown | 40.75362 | 2019-05-21 |
| 2 | Harlem | 40.80902 | NaT |
| 3 | Clinton Hill | 40.68514 | 2019-05-07 |
| 4 | East Harlem | 40.79851 | 2018-11-19 |
| ... | ... | ... | ... |
| 48890 | Bedford-Stuyvesant | 40.67853 | NaT |
| 48891 | Bushwick | 40.70184 | NaT |
| 48892 | Harlem | 40.81475 | NaT |
| 48893 | Hell's Kitchen | 40.75751 | NaT |
| 48894 | Hell's Kitchen | 40.76404 | NaT |

48895 rows × 3 columns

# 5. Missing values

```
# To see the number of missing values
inp0.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
availability_365_categories         0
minimum_night_categories            0
number_of_reviews_categories        0
price_categories                    0
dtype: int64
```

- Two columns (last_review , reviews_per_month) has around 20.56% missing values. name and host_name has 0.3% and 0.4 % missing values
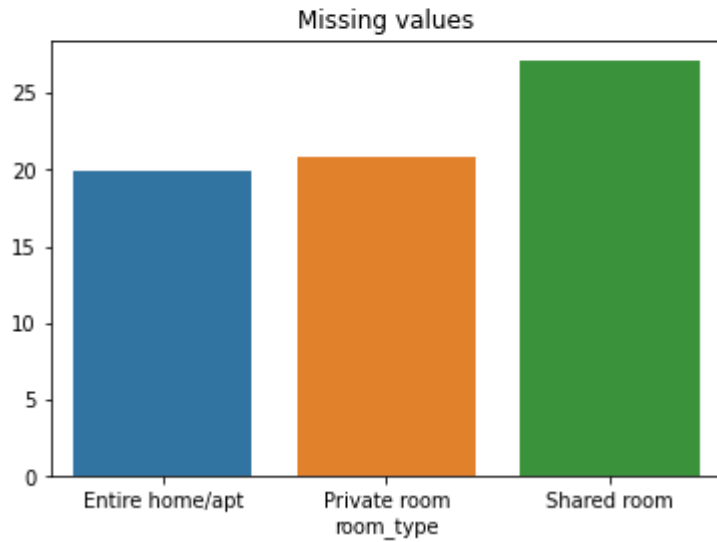
- We need to see if the values are, MCAR: It stands for Missing completely at random.

The reason behind the missing value is not dependent on any other features or if it is MNAR: It stands for Missing not at random. There is a specific reason behind the missing value.
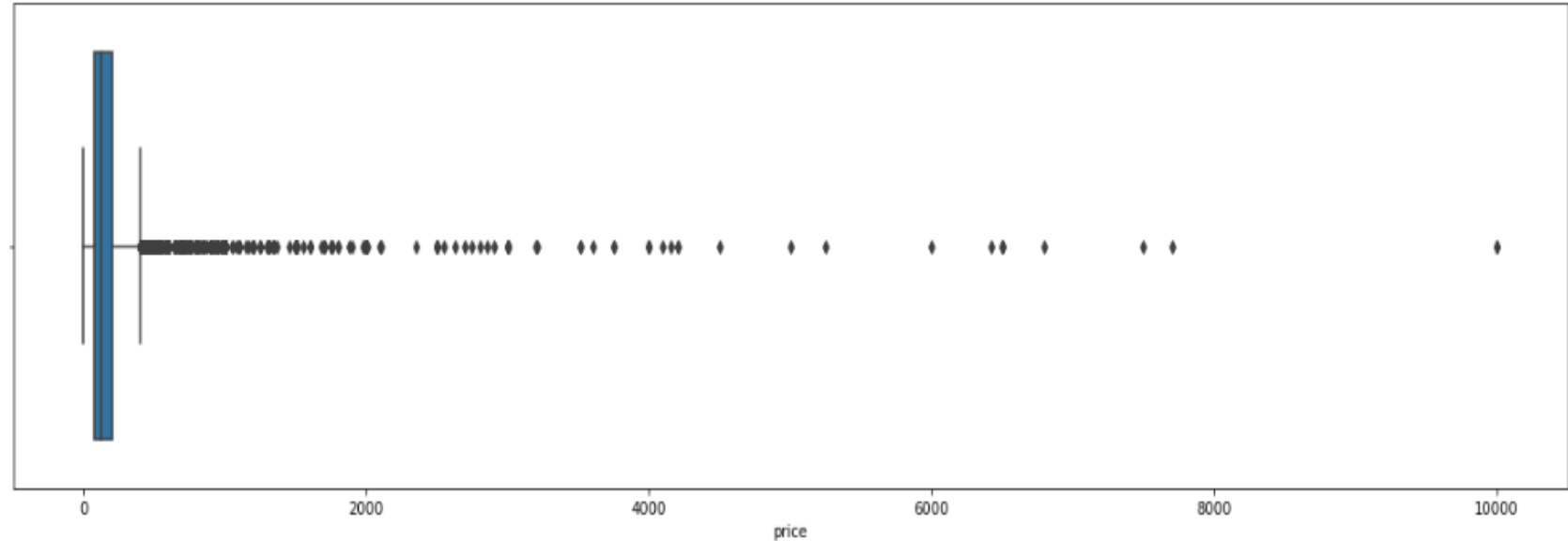
- There is no dropping or imputation of columns as we are just analyzing the dataset and not making a model.Also most of the features are important for our analysis.

# 5.1 Missing value analysis

```
plt.title('Missing values')
sns.barplot(x = inp3.index, y = inp3.values)
plt.show()
```



**'Shared room' has the highest missing value percentage (27 %) for 'last_review' feature while to other room types has only about 20 %.**
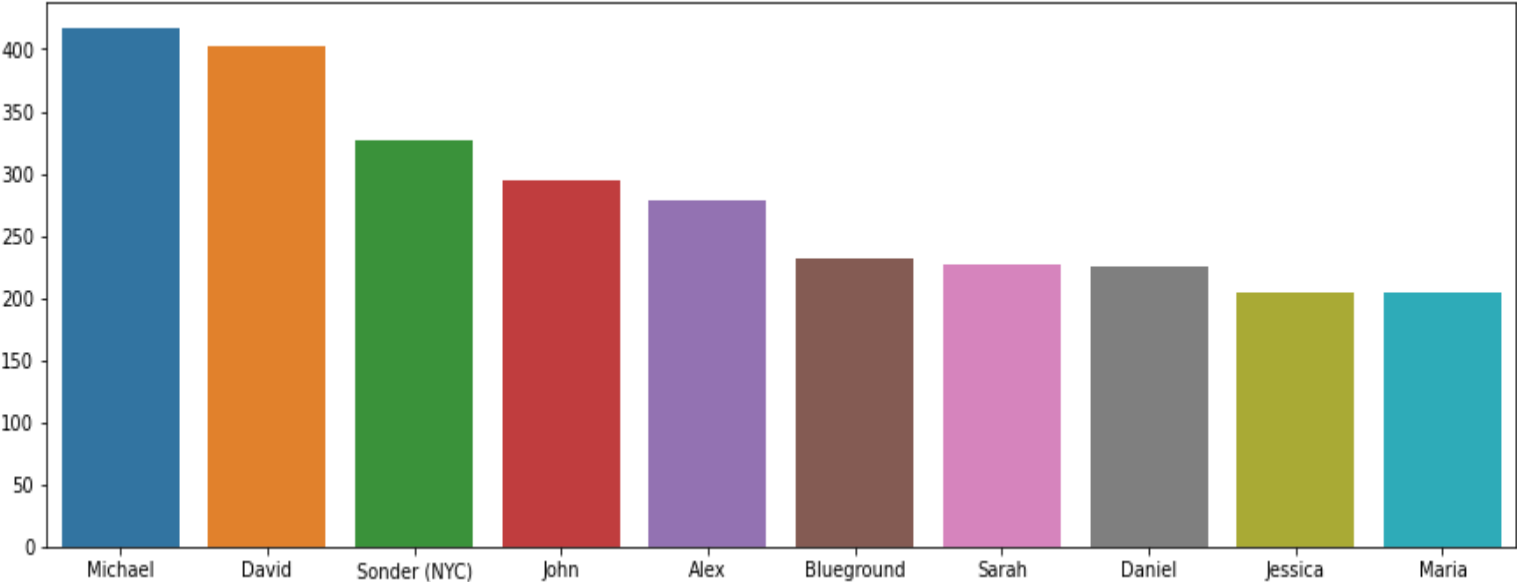


- **The pricing is higher when 'last_review' feature is missing .**

- **reviews are less likely to be given for shared rooms**

- **When the prices are high reviews are less likely to be given**

- **The above analysis seems to show that the missing values here are not MCAR (missing completely at random)**

# 6. Analysis

## 6.3 host_name

```
inp0.host_name.value_counts()
```

```
Michael                  417
David                    403
Sonder (NYC)             327
John                     294
Alex                     279
                        ...
Rhonycs                    1
Brandy-Courtney            1
Shanthony                  1
Aurore And Jamila          1
Ilgar & Aysel              1
Name: host_name, Length: 11452, dtype: int64
```
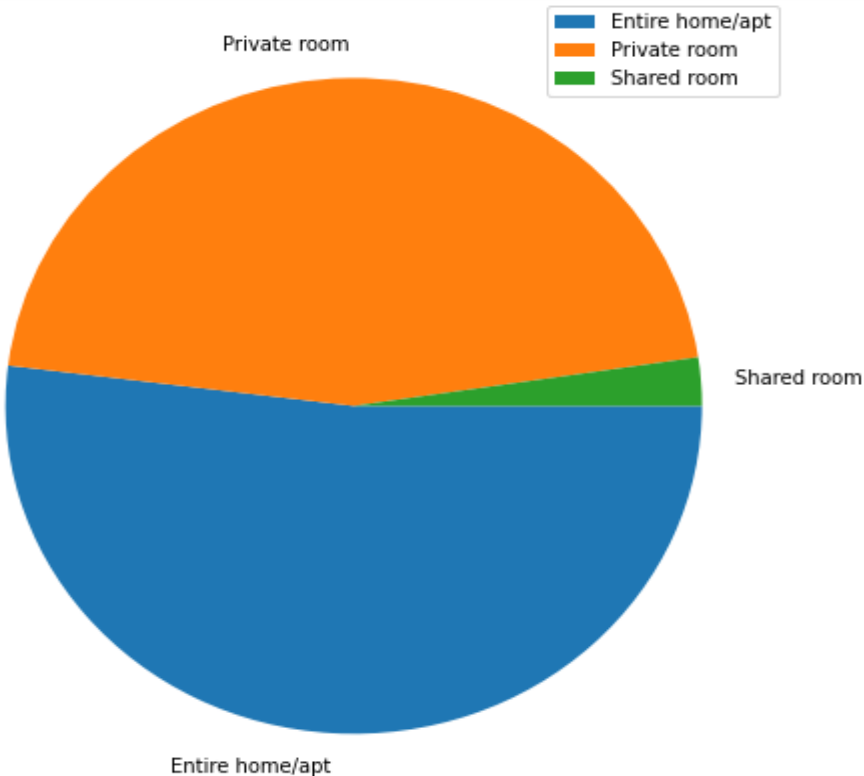


## 6.6 room_type

```
inp0.room_type.value_counts()
```

```
Entire home/apt        25409
Private room           22326
Shared room             1160
Name: room_type, dtype: int64
```
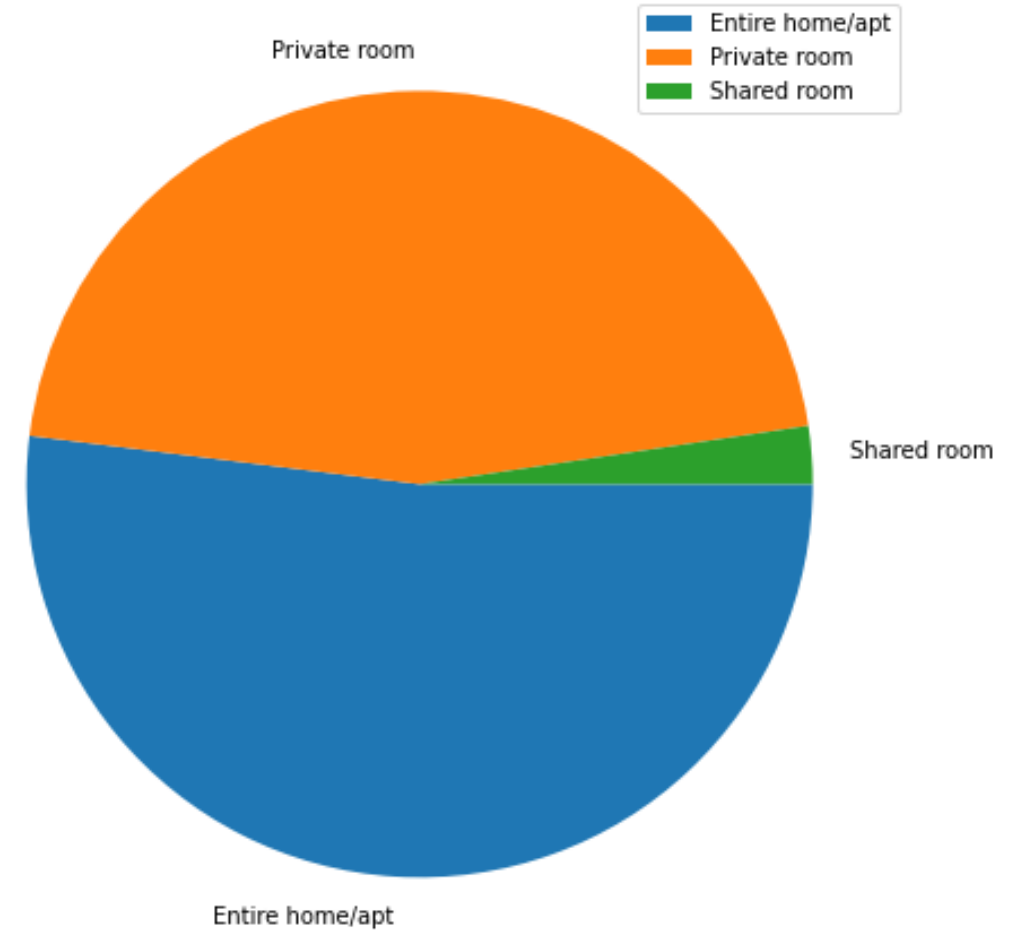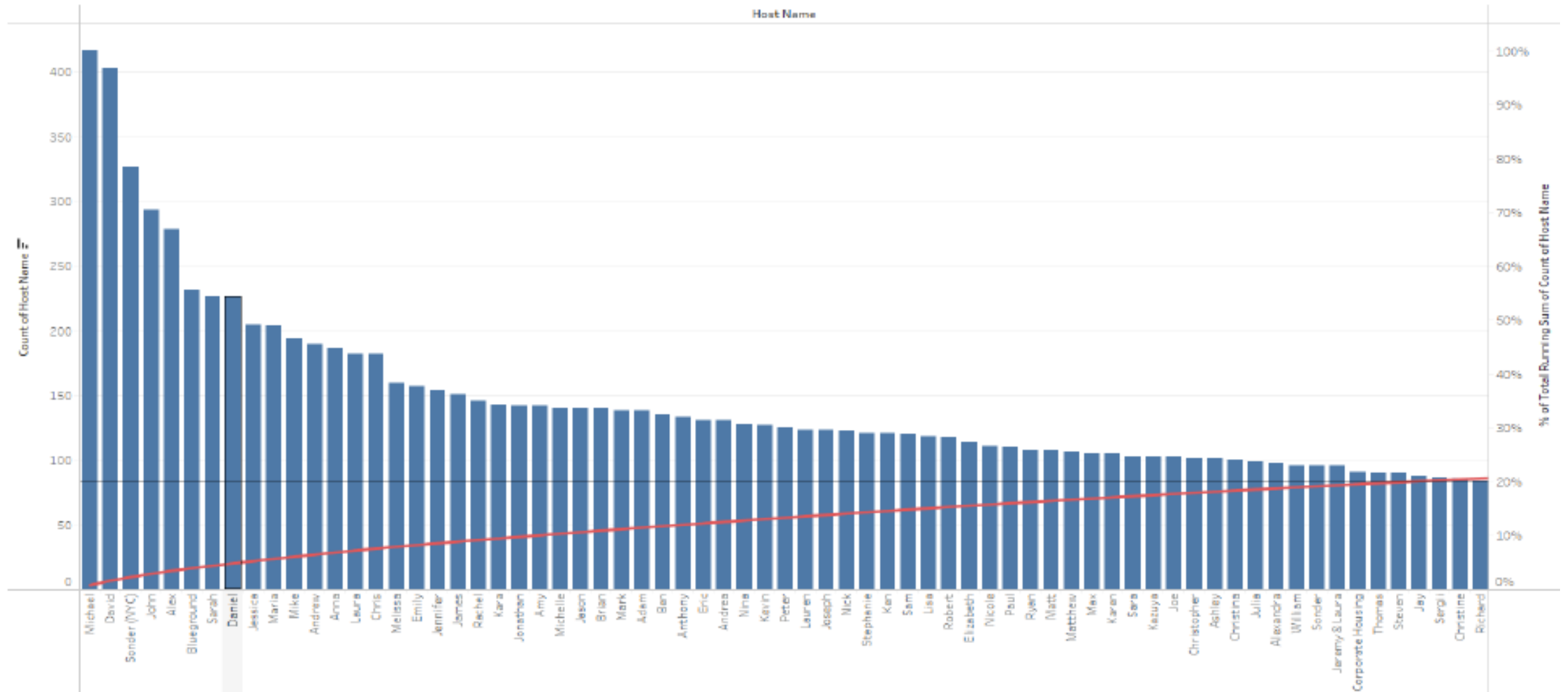
# THE PROBLEMS WITH SHARED ROOMS

**Shared rooms only account for 2 % of the total types of rooms.**

**They are less likely to be reviewed.**

**Median rates for shared rooms are significantly lower.**

# EVERY HOST MATTER



**The top 60 hosts only make up 20% of the total host count!**
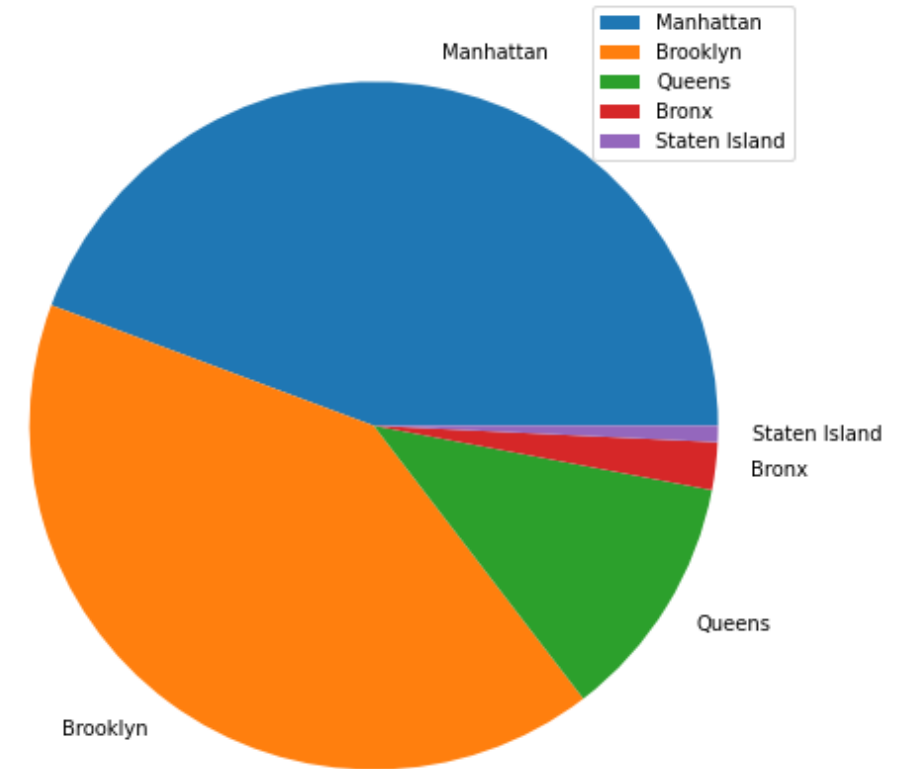
# MOST CONTRIBUTING NEIGHBOURHOODS

**6.4 neighbourhood_group**

```
inp0.neighbourhood_group.value_counts()
```

```
Manhattan          21661
Brooklyn           20104
Queens              5666
Bronx               1091
Staten Island        373
Name: neighbourhood_group, dtype: int64
```
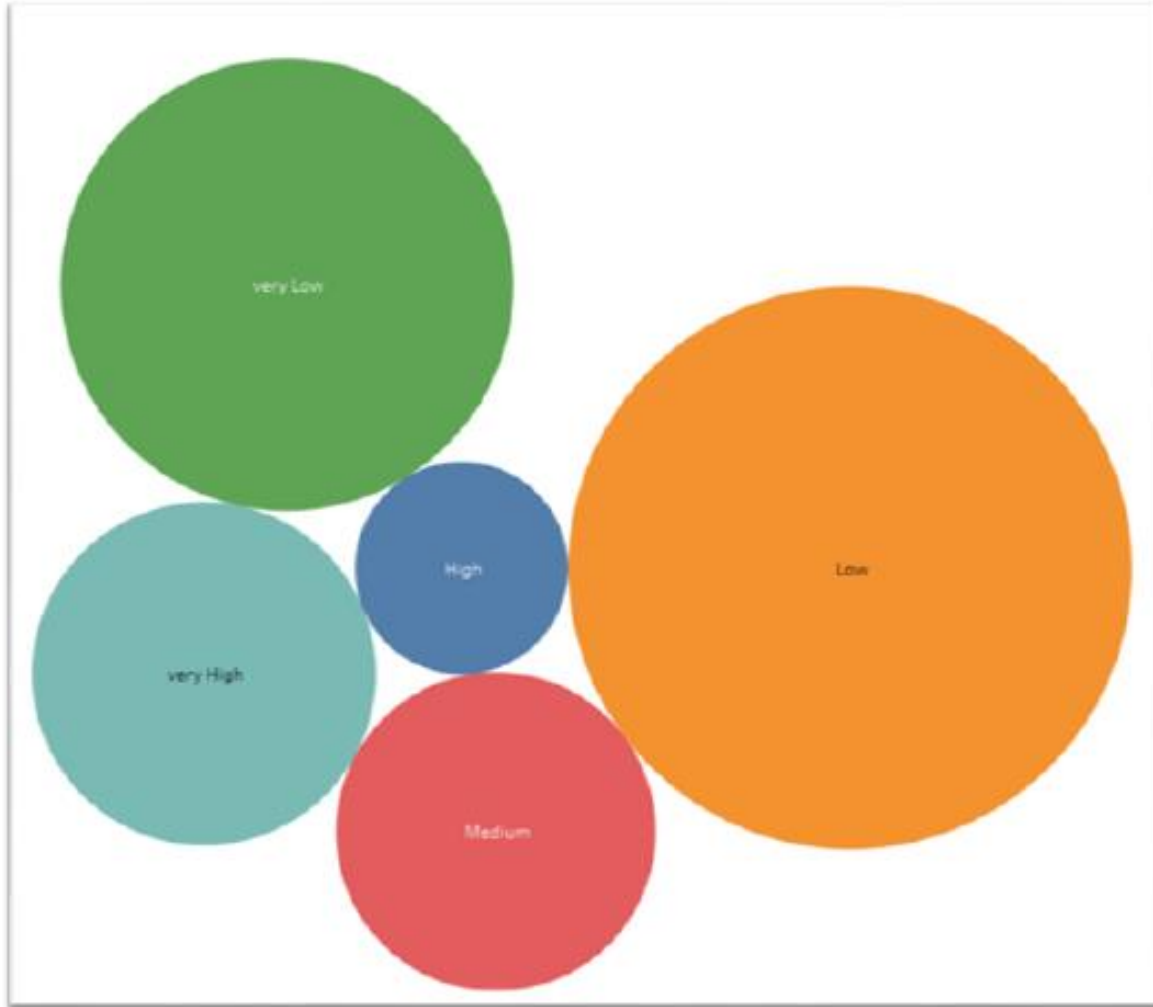


- **81 % of the listing are Manhattan and Brooklyn neighbourhood group**

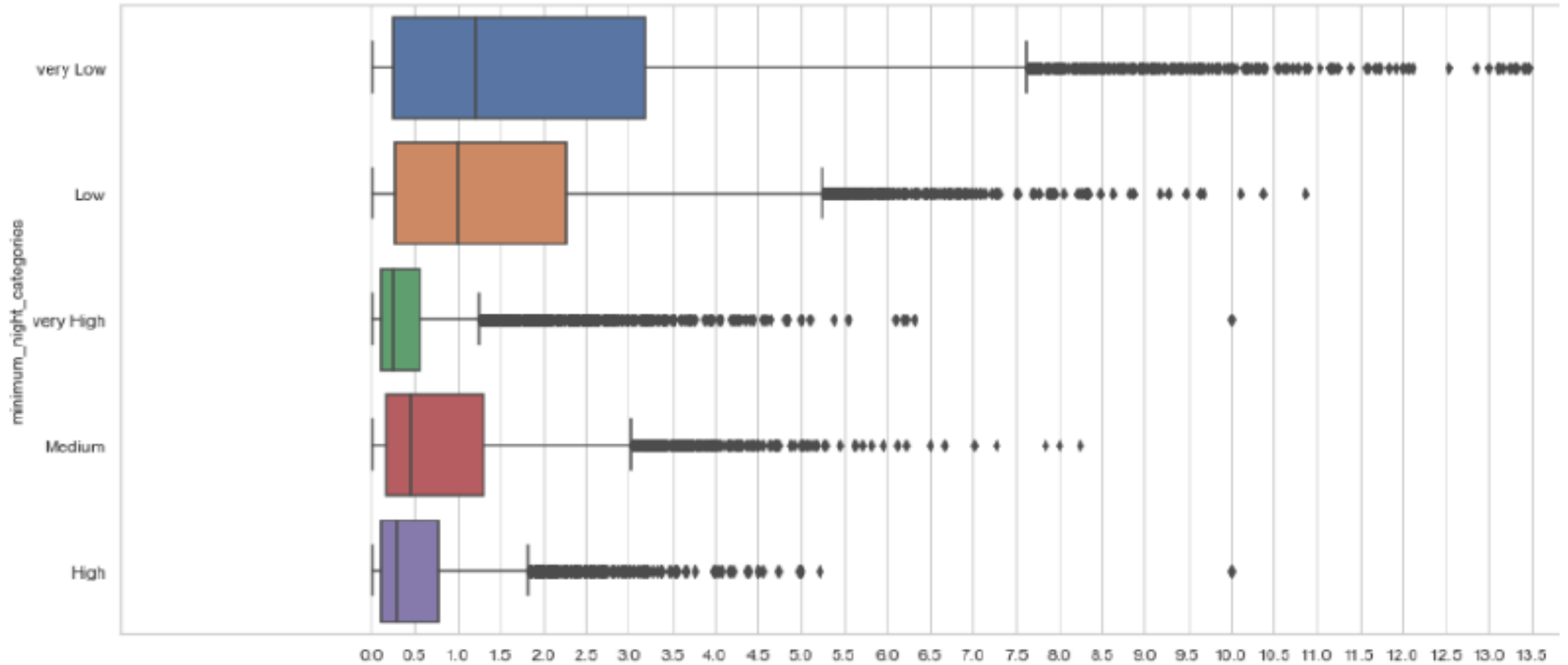- **Staten Island has the lowest contribution.**

# MINIMUM NIGHT CATEGORIES



Minimum night category percentages

| | |
|---|---|
| Low | 40.280192 |
| very Low | 26.014930 |
| very High | 14.997444 |
| Medium | 12.960425 |
| High | 5.747009 |

- **Low category in minimum night feature contributes 40 %**
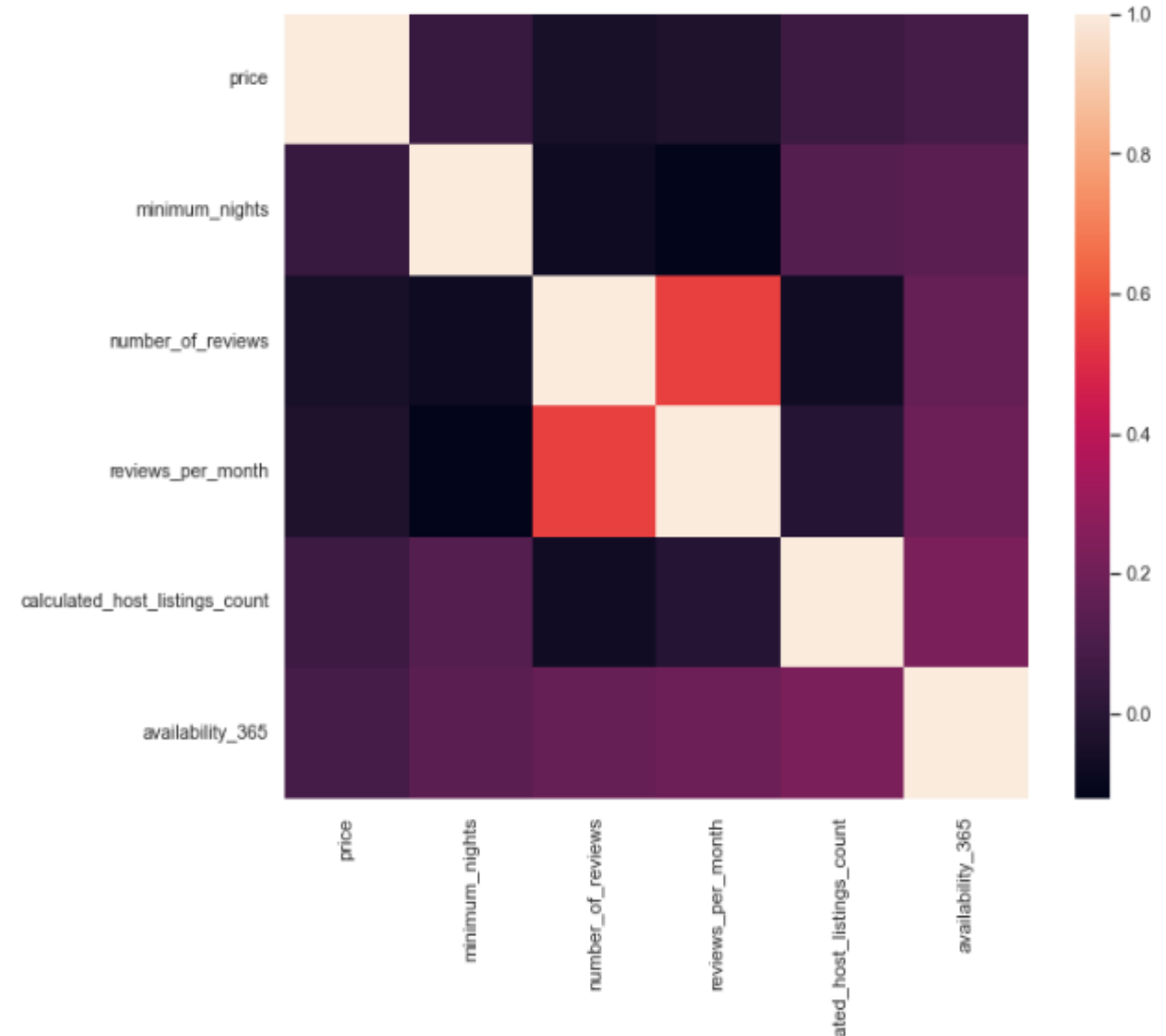
# EFFECT OF MINIMUM NIGHT ON REVIEWS



- **Customers are more likely to leave reviews for lower number of minimum nights.**

# 7. Bivariate and Multivariate Analysis

**7.1 Finding the correalations**

```
inp0[numerical_columns].head()
```

| | price | minimum_nights | number_of_reviews | reviews_per_month | calculated_host_listings_count | availability_365 |
|---|---|---|---|---|---|---|
| 0 | 149 | 1 | 9 | 0.21 | 6 | 365 |
| 1 | 225 | 1 | 45 | 0.38 | 2 | 355 |
| 2 | 150 | 3 | 0 | NaN | 1 | 365 |
| 3 | 89 | 1 | 270 | 4.64 | 1 | 194 |
| 4 | 80 | 10 | 9 | 0.10 | 1 | 0 |

# CONCLUSION

★

1. Strong significant insights are derived based on various attributes in the dataset.

2. Ample amount and variety of visuals have can used in the presentations for the stake-holders.

3. Data collection team should collect data about review scores so that it can strengthen the later analysis.

4. A clustering machine learning model to identify groups of similar objects in datasets with two or more variable quantities can be made.

# APPENDIX -DATA SOURCES

**The columns in the dataset are self-explanatory. You can refer to the diagram given below to get a better idea of what each column signifies.**

| 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 |

# APPENDIX –DATA METHODOLOGY

- **Conducted a thorough analysis of New York Airbnbs Dataset.**

- **Cleaned the data set using python.**

- **Derived the necessary features.**

- **Used group aggregation, pivot table and other statistical methods.**

- **Created charts and visualizations using Tableau.**

# APPENDIX -DATA ASSUMPTIONS

```
Categorical Variables:
     - room_type
     - neighbourhood_group
     - neighbourhood

Continous Variables(Numerical):
     - Price
     - minimum_nights
     - number_of_reviews
     - reviews_per_month
     - calculated_host_listings_count
     - availability_365
- Continous Variables could be binned in to groups too

Location Varibles:
     - latitude
     - longitude

Time Varibale:
     - last_review
```