# STORYTELLING CASE STUDY – AIRBNB NYC

CONTRIBUTORS:

SARAVANAN RAJU

ABHIRAM N

GAYATRI CHAUHAN

# AGENDA

- Objective

- Background

- Key Findings

- Recommendations

- Appendix
  - Data Sources
  - Data Methodology
  - Data Model Assumptions

# OBJECTIVE

- Conduct thorough analysis of Airbnb NYC dataset.

- Ask effective justifications that can lead to data insights.

- Process, Analyse and share findings by Data Visualization and Statistical Techniques.
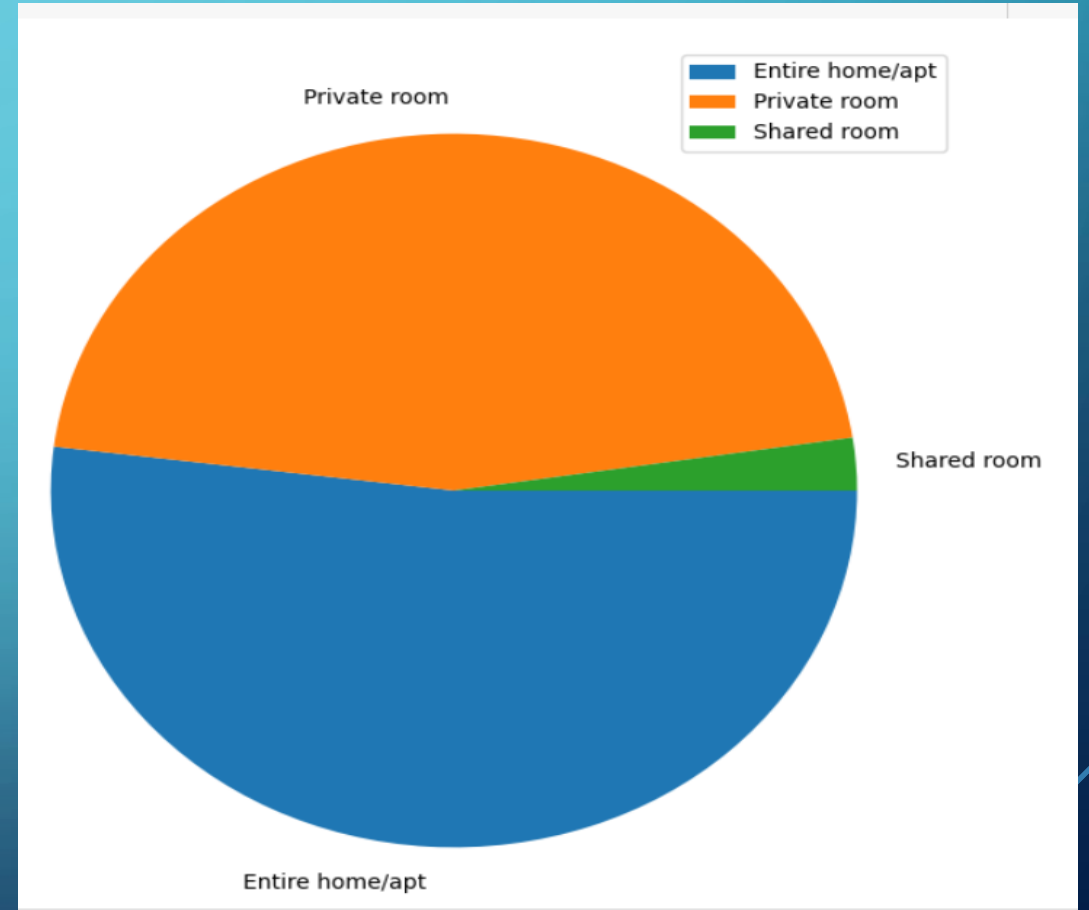
# BACKGROUND

- Airbnb, an online marketplace for short – and long-term homestays and experiences.

- For the past few months, Airbnb has seen a major decline in revenue.

- Now that the restrictions have been lifted and people have started to travel more.

- Airbnb wants to make sure that it is fully prepared for the change.

# DATA ANALYSIS STEPS

- In the first phase the data is captured and loaded for cleansing & data preparation.

- Once data is cleaned, Exploratory data analysis(EDA) is done and new features are created.

- Meaningful insights are created using various analytical methods.

# THE PROBLEM WITH SHARED ROOMS

- Shared rooms only account for 2% of the total type of rooms.
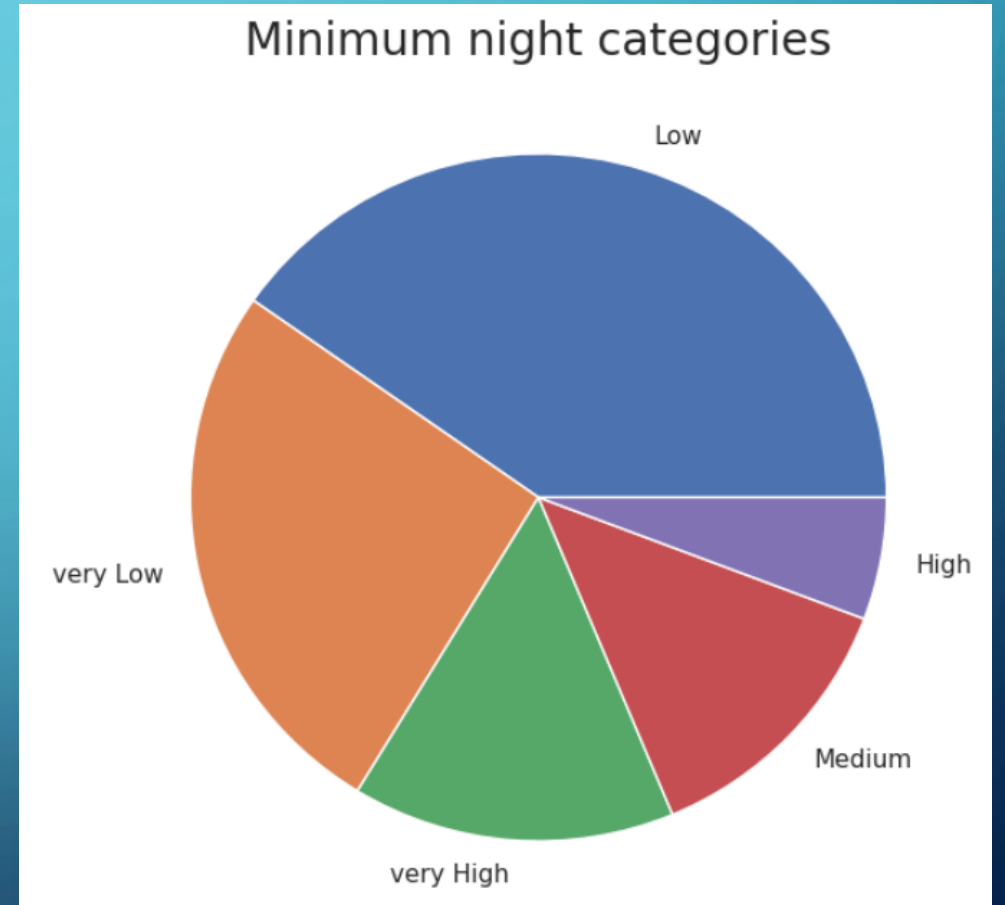
- They are less likely to be reviewed.

# MOST CONTRIBUTING NEIGHBOURHOODS

- 81% of the listings are Manhattan and Brooklyn neighbourhood group.

- Staten Island has the least contribution.

# MINIMUM NIGHT CATEGORIES

- Low category in minimum night feature contributes 40%.



Minimum night categories

# EFFECT OF MINIMUM NIGHTS ON REVIEWS

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



number_of_reviews_categories

# KEY FINDINGS AND RECOMMENDATIONS

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

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

- Shared accommodations has the least preferences. These need to be inspected and customized to private rooms to meet customer demand.

- More than 80% of the listings are in Manhattan and Brooklyn neighbourhood.

- Threshold of minimum nights should be less than 10 nights to make property more customer-oriented.

# APPENDIX

DATA SOURCES

DATA METHODOLOGY

DATA ASSUMPTIONS

# 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.

**Note:** The price column contains the price/night.

| 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 |

**Dataset Description**

# DATA METHODOLOGY

- Conducted Data Analysis on Airbnb, NYC dataset.

- Data-Cleaning, Preparation & adding features were done through Phyton.

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

- Created charts and Visualization through Python & Power-BI

# DATA METHODOLOGY – DATA CLEANING/PREPARATION

# EXPLORATORY DATA ANALYSIS

# EXPLORATORY DATA ANALYSIS – CONTINUED...

# EXPLORATORY DATA ANALYSIS – CONTINUED…

# EXPLORATORY DATA ANALYSIS – CONTINUED...

# EXPLORATORY DATA ANALYSIS – CONTINUED…

# DATA METHODOLOGY – VISUALIZATION THRU POWER BI

# 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
```

**Variable Categories**