

# Capstone Project Submission

## Instructions:

- i) Please fill in all the required information.
- ii) Avoid grammatical errors.

### **Team Member's Name, Email and Contribution:**

**Vivek Kumar Soni:** [vivekkumar75251@gmail.com](mailto:vivekkumar75251@gmail.com)

- 1.Data cleaning:-Handling null values.
- 2.Host\_area,area\_reviews,price\_areas,busiest\_hosts
- 3.Apply dropna () function to entire dataset and drop null values.
- 4.Perform the activity of bar plot using matplotlib of predication and busisest hotels.

**Ashish Kumar pandey:** [ap90920@gmail.com](mailto:ap90920@gmail.com)

- 1.Perform the activity of data cleaning and data wrangling
2. Perform the activity about different hosts and areas
- 3.perform the activity we learn from predictions
- 4.Which hosts are the busiest .
5. Perform noticeable difference of traffic among different areas and what could be the reason for it.
- 6.Plot the bar graph using the matplotlib.

### **Please paste the GitHub Repo link.**

Vivek's Github Link:- <https://github.com/vivek7525/airbnb-data-analytics>

Ashish kumar pandey github link:- <https://github.com/Ashish6681/EDA-AIRBNB-CAPSTONE-PROJECT/tree/main>

## Summary

Ashish kumar pandey

Airbnb guests and hosts have used Airbnb to expand on traveling possibilities and present a more unique, personalized way of experiencing the world. Today, Airbnb became one of a kind service that is used and recognized by the whole world. Data analysis on millions of listings provided through Airbnb is a crucial factor for the company. These millions of listings generate a lot of data - data that can be analyzed and used for security, business decisions, understanding of customers' and providers' (hosts) behavior and performance on the platform, guiding marketing initiatives, implementation of innovative additional services and much more.

## Data Source

This dataset has around 49,000 observations in it with 16 columns and it is a mix between categorical and numeric values.

## Acquaring and Loading Data

After loading the dataset in and from the head of AB\_2019\_NYC dataset we can see a number of things. These 16 columns provide a very rich amount of information for deep data exploration we can do on this dataset. We do already see some missing values, which will require cleaning and handling of NaN values. Later, we may need to continue with mapping certain values to ones and zeros for predictive analytics.

## UNDERSTANDING, DATA ARANGLING AND DATA CLEANING

In our case, missing data that is observed does not need too much special treatment. Looking into the nature of our dataset we can state further things: columns "name" and "host\_name" are irrelevant and insignificant to our data analysis, columns "last\_review" and "review\_per\_month" need very simple handling. To elaborate, "last\_review" is date; if there were no reviews for the listing - date simply will not exist. In our case, this column is irrelevant and insignificant therefore appending those values is not needed. For "review\_per\_month" column we can simply append it with 0.0 for missing values; we can see that in "number\_of\_review" that column will have a 0, therefore following this logic with 0 total reviews there will be 0.0 rate of reviews per month. Therefore, let's proceed with removing columns that are not important and handling of missing data.

Please note that we are dropping 'host\_name' not only because it is insignificant but also for ethical reasons. There should be no reasoning to continue data exploration and model training (which we will be doing later) towards specific individuals based on their names. Why is that? Those names are assigned to actual humans, also they present no security threat or military/governmental interest based on the nature of the dataset, therefore names are unimportant to us.

## Exploring and Visualizing Data

Now that we are ready for an exploration of our data, we can make a rule that we are going to be working from left to right. The reason some may prefer to do this is due to its set approach - some datasets have a big number of attributes, plus this way we will remember to explore each column individually to make sure we learn as much as we can about our dataset.

Amazing, but let's breakdown on what we can see from this plot. First, we can see that our plot consists of 3 subplots - that is the power of using catplot; with such output, we can easily proceed with comparing distributions among interesting attributes. Y and X axes stay exactly the same for each subplot, Y-axis represents a count of observations and X-axis observations we want to count.. So, what do we learn from this? The observation that is definitely contrasted the most is that 'Shared room' type Airbnb listing is barely available among 10 most listing-populated neighborhoods. Then, we can see that for these 10 neighborhoods only 2 boroughs are represented: Manhattan and Brooklyn; that was somewhat expected as Manhattan and Brooklyn are one of the most traveled destinations, therefore would have the most listing availability. We can also observe that Bedford-Stuyvesant and Williamsburg are the most popular for Manhattan borough, and Harlem for Brooklyn.

