CIS 450 Project Proposal - Group 24

1. Team Members

- a. Aakash Jajoo <u>aakashj1@sas.upenn.edu</u> GitHub: aakashjajoo1
- b. Karan Jaisingh karanj@seas.upenn.edu GitHub: kjaisingh
- c. Yathushan Nadanapathan <u>yathu@seas.upenn.edu</u> GitHub: yathu-n
- d. Kush Pandey kpandey@seas.upenn.edu GitHub: kushpandey1811

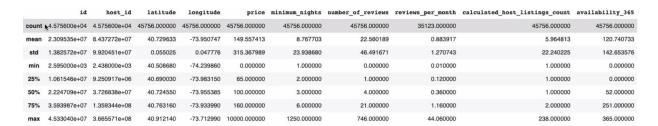
2. Application Description

We are planning on using an airbnb dataset along with dataset on some cities to give recommendations about locations/airbnbs a traveller can stay in depending on their preferences. Additionally, we are going to analyze the correlations between airbnbs (price, rating, reviews, etc) and the number of noise complaints in NYC.

3. Dataset Description

<u>Dataset #1</u>: AirBNB Data - https://www.kaggle.com/kritikseth/us-airbnb-open-data

• This dataset contains information about Airbnb listings in the US. The attributes in the dataset include the name of the host (and the ID of the host in their database), location of the listing (neighborhood and latitude/longitude), type of property and the cost per night. Based on some exploratory analysis, there are 30 unique cities in the dataset and about 226,000 listings. About 45,756 of the listings are in New York City and these are the ones we will focus on, as explained below.



• This is a summary chart of all the listings in the dataset. The data still needs to be cleaned, which is why we see listings with a price of 0 and other noisy data. The number of reviews vary between 0 and 966, while the price varies between 0 (corrupted data) and 24,999 dollars per night. The other statistics can be seen in the table.

| | | | | | | | _ | | | |
|-----------------|---|-------------------|-------------------|----------------|---------------|---------------|---------------|--------------|--------------|-------|
| availability_36 | ${\tt calculated_host_listings_count}$ | reviews_per_month | number_of_reviews | minimum_nights | price | longitude | latitude | host_id | id | |
| 226030.00000 | 226030.000000 | 177428.00000 | 226030.000000 | 2.260300e+05 | 226030.000000 | 226030.000000 | 226030.000000 | 2.260300e+05 | 2.260300e+05 | count |
| 159.31485 | 16.698562 | 1.43145 | 34.506530 | 4.525490e+02 | 219.716529 | -103.220662 | 35.662829 | 9.352385e+07 | 2.547176e+07 | mean |
| 140.17962 | 51.068966 | 1.68321 | 63.602914 | 2.103376e+05 | 570.353609 | 26.222091 | 6.849855 | 9.827422e+07 | 1.317814e+07 | std |
| 0.00000 | 1.000000 | 0.01000 | 0.000000 | 1.000000e+00 | 0.000000 | -159.714900 | 18.920990 | 2.300000e+01 | 1.090000e+02 | min |
| 0.00000 | 1.000000 | 0.23000 | 1.000000 | 1.000000e+00 | 75.000000 | -118.598115 | 32.761783 | 1.399275e+07 | 1.515890e+07 | 25% |
| 140.00000 | 2.000000 | 0.81000 | 8.000000 | 2.000000e+00 | 121.000000 | -97.817200 | 37.261125 | 5.138266e+07 | 2.590916e+07 | 50% |
| 311.00000 | 6.000000 | 2.06000 | 39.000000 | 7.000000e+00 | 201.000000 | -76.919322 | 40.724038 | 1.497179e+08 | 3.772624e+07 | 75% |
| 365.00000 | 593.000000 | 44.06000 | 966.000000 | 1.000000e+08 | 24999.000000 | -70.995950 | 47.734620 | 3.679176e+08 | 4.556085e+07 | max |

This is a summary chart of the listings in New York City. The data still needs to be cleaned,
 which is why we see listings with a price of 0 and other noisy data. The number of reviews vary

between 0 and 746, while the price varies between 0 (corrupted data) and 10,000 dollars per night.

Dataset #2: New York Party Data -

https://www.kaggle.com/somesnm/partynyc?select=bar_locations.csv

- A dataset on the number of noise complaints in the 5 boroughs of NYC between 2015 and 2017.
 Additionally, we have a dataset on 2000 bars in NYC; specifically, we have the bar's latitude and longitude and the number of complaints they've had.
- This dataset contains 225,000 noise complaints (rows). Attributes included in the dataset: city, borough, latitude, longitude, date that the complaint was open, date that the complaint was closed, and location type of the complaint. The average number of complaints that a bar had between 2015 and 2017 is 37.
- Statistics on the party dataset:

| | Incident Zip | Latitude | Longitude | num_calls |
|-------|--------------|-------------|-------------|-------------|
| count | 2440.000000 | 2440.000000 | 2440.000000 | 2440.000000 |
| mean | 10631.856967 | 40.733985 | -73.952497 | 37.025820 |
| std | 591.657847 | 0.066349 | 0.060292 | 59.641884 |
| min | 10001.000000 | 40.511255 | -74.251277 | 10.000000 |
| 25% | 10019.000000 | 40.702064 | -73.987822 | 14.000000 |
| 50% | 10463.000000 | 40.728351 | -73.957928 | 21.000000 |
| 75% | 11217.000000 | 40.765226 | -73.925422 | 40.000000 |
| max | 11694.000000 | 40.910201 | -73.709219 | 1513.000000 |

4. Queries List

<u>Query 1</u>: Average prices and ratings of AirBNBs given a certain geographical location (based on longitude and latitude) or area.

Query 2: Average prices and ratings of AirBNBs given a set of filtering queries, such as room type, host rating and other metrics.

Query 3: Average minimum number of nights a stay needs to be booked for in a particular neighbourhood.

<u>Query 4</u>: Average party activity in a certain geographical location in New York City (based on longitude and latitude) or suburb.

<u>Query 5</u>: Recommended AirBNBs based on the amount of partying that a given individual desires (currently have the data for NYC and plan to add more cities if we can find the relevant data).