

STA 160 Group 7 – Midterm Project

Rishi Bhuva,¹ Aditya Kallepalli,² Vishnu Rangiah,³ and Ivan Yang

¹*Student ID: 915218518, Email: rdhhuva@ucdavis.edu*

²*Student ID: 915079375, Email: arkallepalli@ucdavis.edu*

³*Student ID: 916562849, Email: vrrangiah@ucdavis.edu*

⁴*Student ID: 915463046, Email: igyang@ucdavis.edu*

(we will contact the first author for any question about the article)

In this report, our group performed data analysis on 2019 Airbnb listing data in New York City, USA. Apart from general exploratory data analysis, We specifically performed Natural Language Processing Methods on the listing names to gain further insight on diction in influencing different variables in our dataset. We conclude that hosts with lower-priced listings prioritize appealing towards everyday customers with borough names and positive descriptors, while hosts with more expensive listings prioritize listing property features to appeal to higher-end customers. Additionally, hosts will use different languages to target customers belonging to specific cultural groups.

(95 words)

1 **I. INTRODUCTION**

2 We will be analyzing Airbnb data from the "Airbnb Open New York City 2019 Data Set." The data
3 we are currently analyzing comes from <https://www.kaggle.com/dgomonov/new-york-city-airbnb-open-data>.

5 Airbnb is an online marketplace for lodging, primarily home-stays for vacation rentals, and tourism
6 activities. This platform is accessible through their website or mobile application.

7 This data set describes listing activity and other metrics in New York City, USA in 2019, such as
8 geospatial information, average number of reviews per month, etc. Through this data set, we can get
9 information regarding hosts, geographical availability and necessary metrics needed to make
10 predictions and draw conclusions.

11 We will be performing our analysis using Python with the packages *nltk*, *spacy*, *word cloud* for NLP and
12 *pandas*, *matplotlib*, *seaborn*, *numpy*, *geopandas*, *contextily* for EDA. In Section 2, we will be cleaning and
13 performing exploratory analysis on the data set. Then we will explore the relationship between the
14 names of each listing and the popularity and price of each listing using Natural Language Processing
15 Methods in Section 4. We wanted to specifically look at word choice in listing names as they are
16 usually the first piece of information that a customer encounters when searching for listings on the
17 Airbnb platform. Interestingly, we observed a wide variation, ranging from language choice, image-
18 provoking adjectives, to even the use of emojis. Examples of headings:

19 **II. EXPLORATORY DATA ANALYSIS**

20 **A. Main Data Frame**

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
0	Clean & quiet apt home by the park	2787	John	Brooklyn	Kensington	40.64749	-73.97237	Private room	149	1	9	2018-10-19	0.21	6	365
1	Skyit Midtown Castle	2845	Jennifer	Manhattan	Midtown	40.75368	-73.98377	Entire home/apt	225	1	45	2019-05-21	0.38	2	355
2	THE VILLAGE OF HARLEM...NEW YORK!	4632	Elisabeth	Manhattan	Harlem	40.80906	-73.94190	Private room	150	3	0	Nan	Nan	1	365
3	Cozy Entire Floor of Brownstone	4869	LisaRoxanne	Brooklyn	Clinton Hill	40.68514	-73.95976	Entire home/apt	89	1	270	2019-07-05	4.64	1	194
4	Entire Apt: Spacious StudioLoft by central park	7192	Laura	Manhattan	East Harlem	40.79851	-73.94399	Entire home/apt	80	10	9	2018-11-19	0.10	1	0

21 FIG. 1. (Main Data Frame).

22 Here we find that our dataset contains 48,895 rows which represent individual listings in New
 23 York. We are also provided with 16 different columns which contain necessary factors needed for
 24 exploratory data analysis. They are: id, name, host_id, host_name, neighbourhood_group,
 25 neighbourhood, latitude, longitude, room_type, price, minimum_nights, number_of_reviews,
 26 last_review, reviews_per_month, calculated_host_listings_count, and availability_365. 7 variables are
 27 of integer type, 6 variables are floats, and the others are objects.

29 **B. Data Cleaning**

30 Our first line of action when analyzing this data is to properly clean the dataset. Having clean data
 31 will provide us with the highest quality of information needed and therefore will provide us with the
 32 most accurate predictions and correlations. Our process for cleaning the dataset can be seen below.

33 1. Removing Missing and Irrelevant Information

34 We will first check the number of rows with missing values. We will then remove such rows
 35 depending on the missing variables. There are 4 variables with missing values. However, we've
 36 determined that all of these missing variables do not greatly affect our NLP analysis. Therefore,
 37 we've decided to keep these listings.

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
4	Entire Apt: Spacious StudioLoft by central park	7192	Laura	Manhattan	East Harlem	40.79851	-73.94399	Entire home/apt	80	10	9	2018-11-19	0.10	1	0
6	BlissArtsSpace!	7356	Garon	Brooklyn	Bedford-Stuyvesant	40.68688	-73.95598	Private room	60	45	49	2017-10-05	0.40	1	0
8	Cozy Clean Guest Room - Family Apt	7490	MaryEllen	Manhattan	Upper West Side	40.80178	-73.98723	Private room	79	2	118	2017-07-21	0.99	1	0
14	West Village Nest - Superhost	11975	Aina	Manhattan	West Village	40.75530	-74.00525	Entire home/apt	120	90	27	2018-10-31	0.22	1	0
20	Sweet and Spacious Brooklyn Loft	21207	Chaya	Brooklyn	Williamsburg	40.71942	-73.95718	Entire home/apt	299	3	9	2011-12-28	0.07	1	0

39 FIG. 2. (Missing Values in Availability Feature of Data Frame).

40 On the other hand, while combing through our main dataframe, we discovered that there are a
41 number of listings with 0 available days to be rented as well as 0 total reviews. We chose to extract
42 these listings from the dataframe we will analyze, because zero availability means that the listings are
43 not available at all. Therefore, this data would be irrelevant as we want to perform analysis on
44 listings that have availability.

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	
2	3847	THE VILLAGE OF HARLEM...NEW YORK !	4632	Elisabeth	Manhattan	Harlem	40.80902	-73.94190	Private room	150	3	0	NaN	NaN	1	365
19	7750	Huge 2 BR Upper East Central Park	17985	Sing	Manhattan	East Harlem	40.79685	-73.94872	Entire home/apt	190	7	0	NaN	NaN	2	249
26	8700	Magnifique Suite au N de Manhattan - vue Cloîtres	26394	Claude & Sophie	Manhattan	Inwood	40.86754	-73.92639	Private room	80	4	0	NaN	NaN	1	0
36	11452	Clean and Quiet in Brooklyn	7355	Vt	Brooklyn	Bedford-Stuyvesant	40.68876	-73.94312	Private room	35	60	0	NaN	NaN	1	365
38	11943	Country space in the city	45445	Harriet	Brooklyn	Flatbush	40.63702	-73.96327	Private room	150	1	0	NaN	NaN	1	365

45 FIG. 3. (Missing Values in Number of Reviews Feature of Data Frame).

46 The data frame above consists of listings with zero number of reviews. We chose to extract
47 these values from the data we want to analyze, because number of reviews is a factor that can tell us
48 a lot about the popularity of the listing and is an important component to our Natural Language
49 Processing Method which we will perform. Therefore, listings with zero number of reviews would
50 be irrelevant data and not needed for our type of analysis.

52 2. Examining Outliers – Interquartile Test (IQR)

53 Now, we will examine the price range of all Airbnb listings by separating these listings into quartiles
54 to get a general insight for this column alone. Along with the 25, 50, 75 % quartile calculations, we
55 can also get information such as the mean, standard deviation and the minimum and maximum
56 value of our price category.

	price
count	26155.0
mean	149.8851844771554
std	198.81696180412806
min	0.0
25%	70.0
50%	109.0
75%	175.0
max	9999.0

57

58 FIG. 4. (Summary Statistics of Price Feature).

59 As observed in our quartile table above, there are surprisingly listings that are \$0. Our next course of
60 action would be to pull up these listings and do further research.

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
23161	Huge Brooklyn Brownstone Living, Close to it all...	8995084	Kimberly	Brooklyn	Bedford-Stuyvesant	40.69023	-73.95428	Private room	0	4	1	2018-01-06	0.05	4	28
25433	★Hostel Style Room Ideal Traveling Buddies★	131697576	Anisha	Bronx	East Morrisania	40.83296	-73.88668	Private room	0	2	55	2019-06-24	2.56	4	127
25634	MARTIAL LOFT S. REDEMPTION (upscale 2nd room)	15787004	Martial Loft	Brooklyn	Bushwick	40.69467	-73.92433	Private room	0	2	16	2019-05-18	0.71	5	0
25753	Sunny, Quiet Room in Greenpoint	1641537	Lauren	Brooklyn	Greenpoint	40.72462	-73.94072	Private room	0	2	12	2017-10-27	0.53	2	0
25778	Modern apartment in the heart of Williamsburg	10132168	Aymaric	Brooklyn	Williamsburg	40.70838	-73.94645	Entire home/apt	0	5	3	2018-01-02	0.15	1	73
25794	Spacious comfortable master bedroom with nice ...	86327101	Adeyemi	Brooklyn	Bedford-Stuyvesant	40.68173	-73.91342	Private room	0	1	93	2019-06-15	4.28	6	176
25795	Contemporary bedroom in brownstone with nice view	86327101	Adeyemi	Brooklyn	Bedford-Stuyvesant	40.68279	-73.91170	Private room	0	1	95	2019-06-21	4.37	6	232
25796	Cozy yet spacious private brownstone bedroom	86327101	Adeyemi	Brooklyn	Bedford-Stuyvesant	40.68258	-73.91284	Private room	0	1	95	2019-06-23	4.35	6	222
26259	the best you can find	13709292	Quchi	Manhattan	Murray Hill	40.75091	-73.97597	Entire home/apt	0	3	0	NaN	NaN	1	0
26841	Coliving in Brooklyn! Modern design / Shared room	101970559	Sergi	Brooklyn	Bushwick	40.69211	-73.90567	Shared room	0	30	2	2019-06-22	0.11	6	333
26866	Best Coliving space ever! Shared room.	101970559	Sergi	Brooklyn	Bushwick	40.69166	-73.90928	Shared room	0	30	5	2019-05-24	0.26	6	139

61

62 FIG. 5. (Low Values in Price Feature of Data Frame).

63 After multiple online searches of these listings, we've come up with mixed results of these listings'
64 robustness. Listings such as Kimberly's seem to have been removed while all three of Adeyemi's
65 listings are still up on Airbnb and running in 2021. Additionally, there isn't any way that we can
66 revisit Airbnb's 2019 websites to verify these listings.

67 One possible reason for these \$0 prices may have been a deliberate act by the hosts to temporarily
68 remove the listing from the Airbnb market when these listings were webscraped. The additional fact
69 that among these \$0 listings, 3 and 2 of them belong to the same host may further indicate that it is a
70 host-activated anomaly.

71 Either way, we will remove these listings just to minimize any possibility of accruing errors based on
72 unknown anomalies. Below are our summary statistics with entries of price = \$0 excluded.

	price
count	26147.0
mean	149.93104371438406
std	198.83008554328674
min	10.0
25%	70.0
50%	109.0
75%	175.0
max	9999.0

73

74 FIG. 6. (Summary Statistics of Price Feature

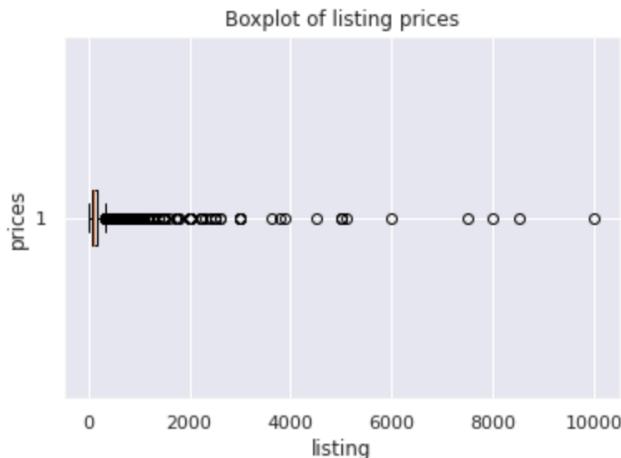


FIG. 7. (Box plot of Price Feature).

75 After Cleaning).

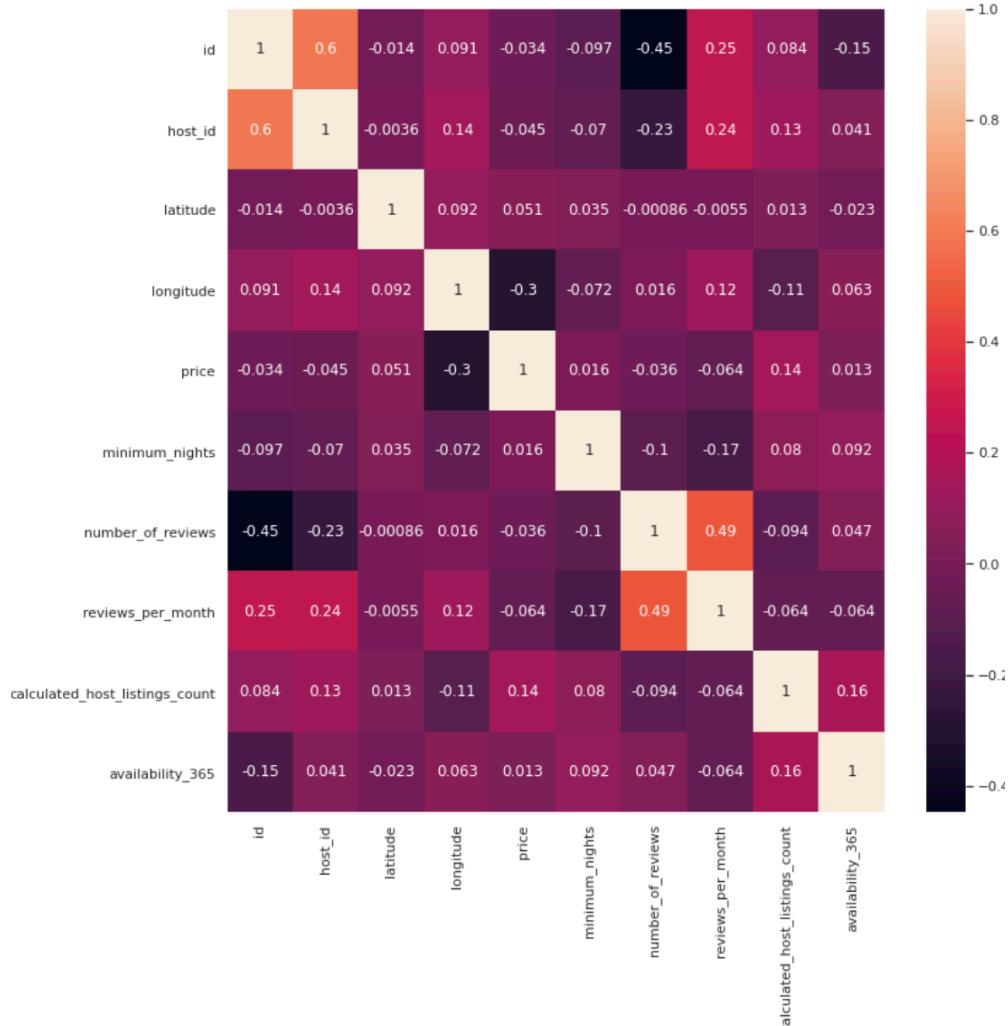
76 From the boxplot regarding our modified prices above, we can see that the data is quite skewed with
77 multiple extremely high prices heavily affecting the data. Further research also shows that some of
78 these data points are faulty listings, such as a 2012 Superbowl private room listing.

79 We will now use the 1.5 Interquartile Range Test to determine outliers within price and store these
80 listings in a different dataframe for further data analysis.

81 Through the 1.5 Interquartile Range Test, we can properly establish the range of outliers. This test
82 essentially says that a data point is considered an outlier if it is more than $1.5 * \text{IQR}$ above the third
83 and first quartile. Therefore low outliers are below our $25^{\text{th}} \text{ percentile} - 1.5 * \text{IQR}$ and high outliers
84 are above our $75^{\text{th}} \text{ percentile} + 1.5 * \text{IQR}$. The IQR variable is calculated using the difference
85 between the 75^{th} percentile and the 25^{th} percentile. The result from this test to determine outliers
86 is: $(-87.5, 332.5)$. We can see that anything above the price of 332.5 dollars would be considered too
87 expensive, and therefore would be an outlier. Since there technically cannot be any negative values
88 for prices, we will fix the lower limit at 10 dollars which is the minimum value as seen in our quartile
89 table. Therefore, any prices under 10 dollars would be considered low outliers. Our next course of
90 action would be to remove these price outliers from the dataset we want to analyze.

91 3. Tidied Dataset

92 After removing all unnecessary information and outliers, we are left with a tidied dataset consisting
 93 of 24540 listings.



94

95 FIG. 8. (Correlation Matrix Between Variables).

96 Here, we can see a proper correlation matrix between the variables within our dataset. According to
 97 this plot, we can see that the more lighter the color, the higher the correlation between each variable.
 98 Vice versa, we can see that the darker the number, the lower the correlation. From this plot, we can
 99 immediately see a high correlation of 0.6 between id and host_id. Along with this, it can be seen that
 100 number of reviews and reviews per month also have a higher correlation than others at 0.49. Along

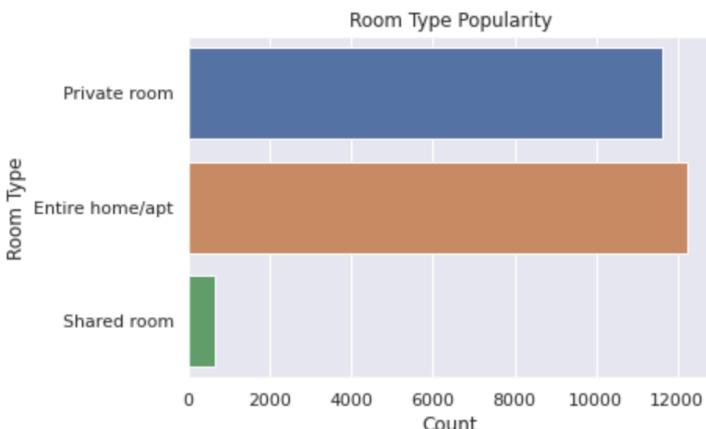
101 with positive correlations, we can also see negative correlations with id and number of reviews at -
102 0.45 and with id and reviews per month at -0.45. We can also see a negative correlation between
103 price and longitude at -0.3, although a positive correlation of 0.051 between price and latitude.
104 As id is a specific tag per listing, host_id is also a specific tag per host. Therefore, this positive
105 correlation does seem to make sense as some hosts would have multiple listings under their name.
106 Number of reviews and reviews per month also tend to have a positive correlation as each number
107 of reviews contributes to a review per month and vice versa, therefore these two variables seem to
108 have a direct relationship.

109 The lowest correlation we can see in this plot can be seen between both variables of reviews and id.
110 The large negative correlation between them does seem accurate, as each id is specific and randomly
111 generated and does not relate at all to reviews at all. Another negative correlation that was interesting
112 comes from price and longitude. As longitude tends to be west and east, it can be seen that price in
113 New York does not tend to have any sort of relationship with those specific directions. Although,
114 price and latitude show a somewhat positive correlation, meaning that price is more correlated with
115 listings that range from the north and south rather than the west and east.

116 III. GRAPHS

117 We will plot some graphs with our tidied dataset to gain further insight, as well as to visualize any
118 more abnormal data points for further cleaning. We will not be subsetting these outliers from the
119 main tidied dataset as these variables will not be our main focus for our NLP analysis. However, we
120 will generate graphs of these outliers for further analysis.

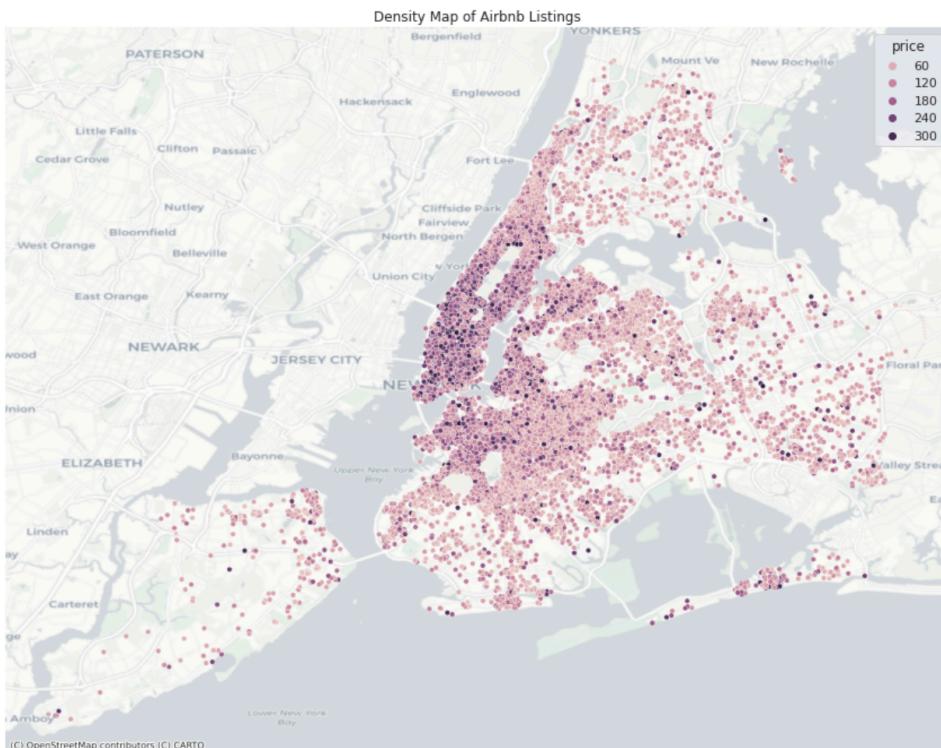
121 A. Tidied Dataset



122

123 FIG. 9. (Bar Plot of Room Type Popularity).

124 We can see that most listings are either private rooms or the entirety of the property. This indicates
125 that most Airbnb customers may be looking for private accommodation rather than shared spaces.

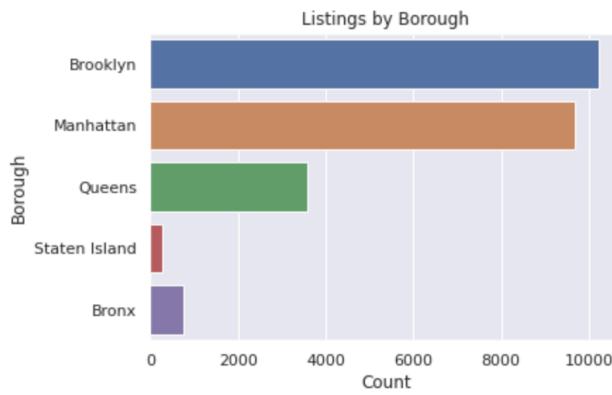


126

127 FIG. 10. (Density Map of Price of Airbnb Listings).

128 From the density map above, we can see that most of the listings are concentrated around the
129 boroughs of Brooklyn and Manhattan. We can also observe that Manhattan has a higher density of

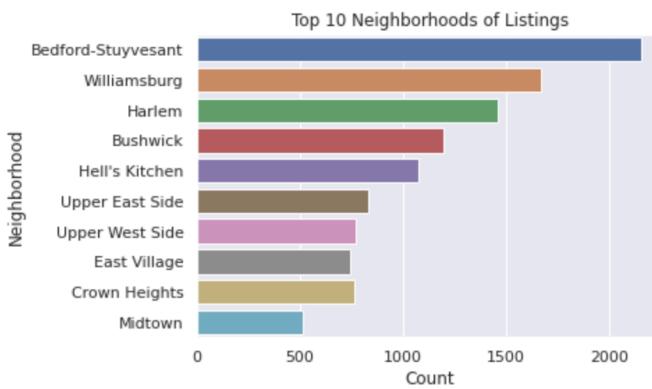
130 expensive listings, and apart from standout anomalies, prices tend to be lower as listings become
131 further from Manhattan.



132

133 FIG. 11. (Bar Plot of Listings by Borough).

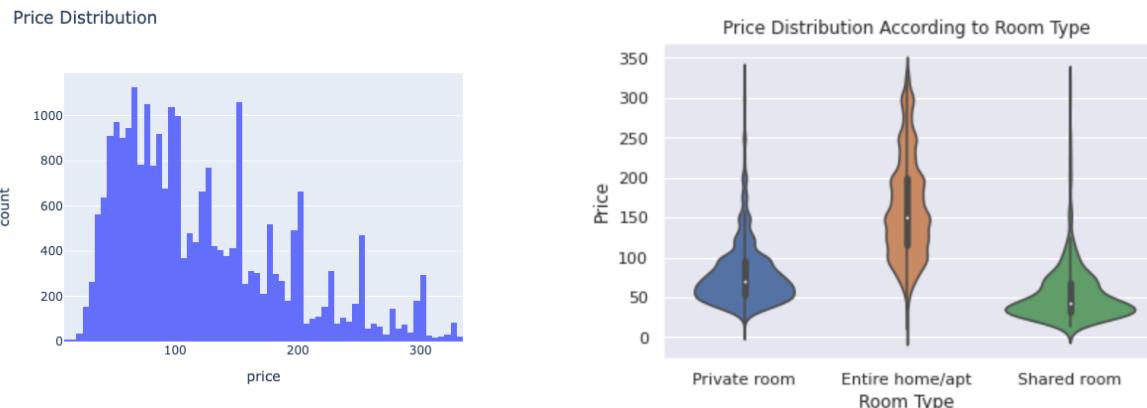
134 A deeper analysis shows that most listings are either in the boroughs of Brooklyn or Manhattan.
135 This is logical given that these 2 boroughs contain the most attractions and activities in New York
136 City. Staten Island has the lowest number of listings, which may be due to its geographical location
137 being relatively inaccessible compared to the other boroughs.



138

139 FIG. 12. (Bar Plot of Listings by Neighborhoods).

140 The top 4 neighborhoods of Airbnb listings are all in Brooklyn, with Bedford-Stuyvesant being the
141 top neighborhood by quite a margin. The 5th to 8th most popular neighborhoods are in Manhattan.
142 This may be due to affordability of Brooklyn compared to Manhattan.



143

144 FIG. 13. (Histogram of Listing Prices).

FIG. 14. (Violin Plot of Room Type).

145 From the histogram above, we can see a left-skewed unimodal distribution. A majority of the
 146 listings hover around a price between 35 to 130 USD. However, there are significant spikes at 150,
 147 200, 250, and 300 USD bins, as well as smaller spikes at 175, 225, 275, and 325 USD bins. We
 148 believe that this may be caused by human psychology where hosts may round their prices to the
 149 nearest 25 or 50 dollar for a more "aesthetically pleasing" price number.

150 The violin plot above shows the distribution of price by room type. There is much variation
 151 in price within each room type. We can see that price greatly fluctuates between the type of room. It
 152 can be seen that entire home/apt listings tend to be more expensive and evenly distributed across
 153 the price range than both private room and shared room, although private room is slightly more
 154 expensive than a shared room. A majority of the private and shared rooms tend to be around \$40 to
 155 \$60. The logic behind the prices seem to make sense.



156

157 FIG. 15. (Listing Prices Based on Boroughs and Room Types).

158 As also observed in our previous violin plot, there is an obvious trend that shared rooms are the
 159 cheapest option in all boroughs, followed by private room and entire property. Additionally, Staten
 160 Island and the Bronx have a lot less listings above \$150 compared to Brooklyn and Manhattan.
 161 Additionally, as seen in our density map, Manhattan has the highest density of expensive listings
 162 compared to other boroughs.

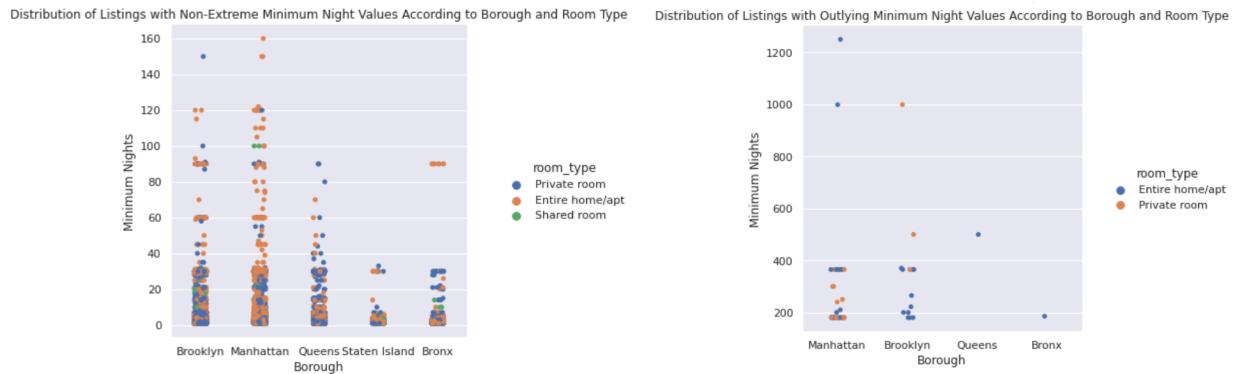
Minimum Nights Distribution



163

164 FIG. 16. (Histogram of Minimum Nights).

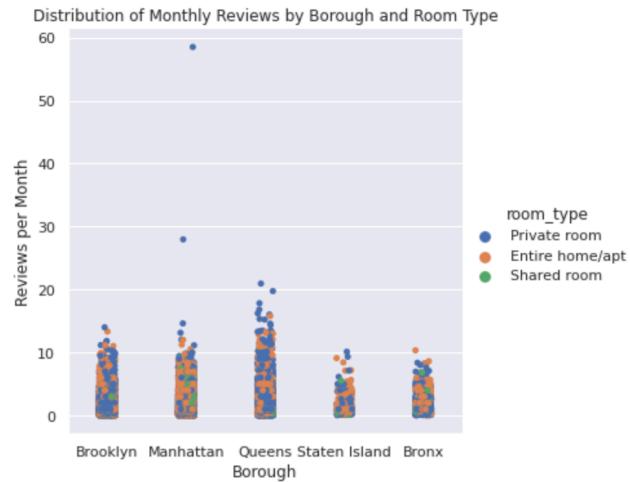
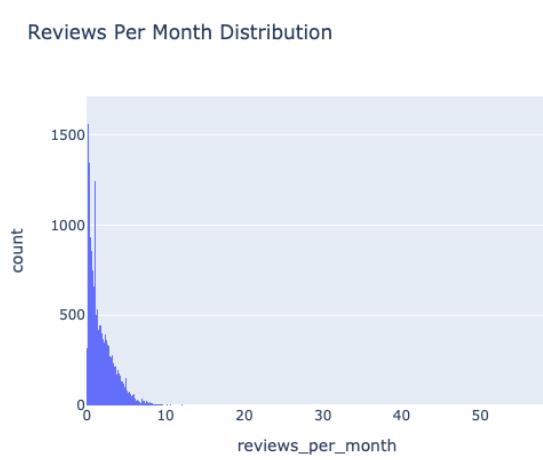
165 Most listings in 2019 NYC have a minimum night requirement between 1 to 7 nights, with an
 166 average minimum night requirement of 6.64 and 2 nights being the highest. There is an apparent
 167 spike of listings with a minimum night requirement of 30 nights.
 168 There seems to be multiple outliers, with one as extreme as 1250 minimum nights. This wide range
 169 of minimum night requirements allude to the existence of not just long term vacation rental
 170 properties but also potentially long term residence.
 171 We will determine these outliers again using the 1.5 IQR test. We obtain a lower limit of -156.5 and
 172 161.5 nights. Since there technically be any negative values for minimum nights, we will fix the lower
 173 limit at 0. We end up with 47 outlying listings with extreme values of minimum nights.



174 FIG. 17. (Minimum Nights based on Borough
 175 and Room Types).
 176 FIG. 18. (Outlying Minimum Nights based
 177 on Borough and Room Type.)
 178 From the graph, it seems like Staten Island and the Bronx are not popular boroughs for long term
 179 rentals longer than 30 days. Further investigation of the few orange anomalies in the Bronx reveal
 180 that they are all from the same host, Sasha (Host ID: 2988712), who seems to have multiple 90-day
 181 Most longer-term rentals in Brooklyn, Manhattan, and Queens are of entire properties. Interesting
 182 anomalies include the 90-day minimum night requirement listings for a shared room in Manhattan.
 183 Further investigation shows that these listings (Host IDs: 21628183 and 23184420) are both from

184 LaGuardia Houses Public Housing Development looking for long term housemates. We believe that
185 these listings are a result of hosts trying to bypass New York's legal barrier of subletting public
186 housing for additional income, adding onto the high accessibility of Airbnb for anyone to put out a
187 listing.

188 As expected, most extreme values of minimum night requirement are from listings in Manhattan and
189 Brooklyn, with Queens and the Bronx each only having 1 listing and Staten Island not having any.
190 Both listings in Queens and the Bronx seem to be long term rental leasings, just like most of the
191 other leasings in this plot.



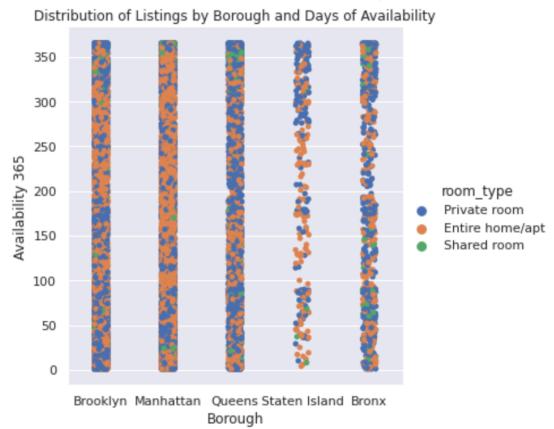
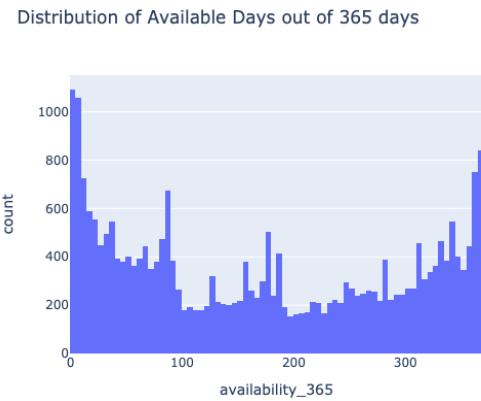
192
193 FIG. 19. (Histogram of Price Feature).

FIG. 20. (Reviews per Month based on Borough and
194 Room Type).

195 The number of reviews are heavily skewed to the left with listings getting an average of 1.83
196 reviews per month. There are obvious outliers, with a listing getting as many as 58.5 average reviews
197 per month.

198 We can see that most boroughs have a similar pattern of reviews per month, with 2
199 anomalies in Manhattan. Further analysis reveals that both listings (ID: 32678719 and 32678720) are
200 by the same host Row NYC (Host ID: 244361589) which is a hotel in New York City. They may
201 have more reviews as they may be more attractive to customers as an established hotel chain.

202 Additionally, they may have multiple rooms available under the same listing, which would explain
203 why they have listings getting 58.5 reviews in a span of 30 or 31 days.



204

205 FIG. 21. (Histogram of Available Days).

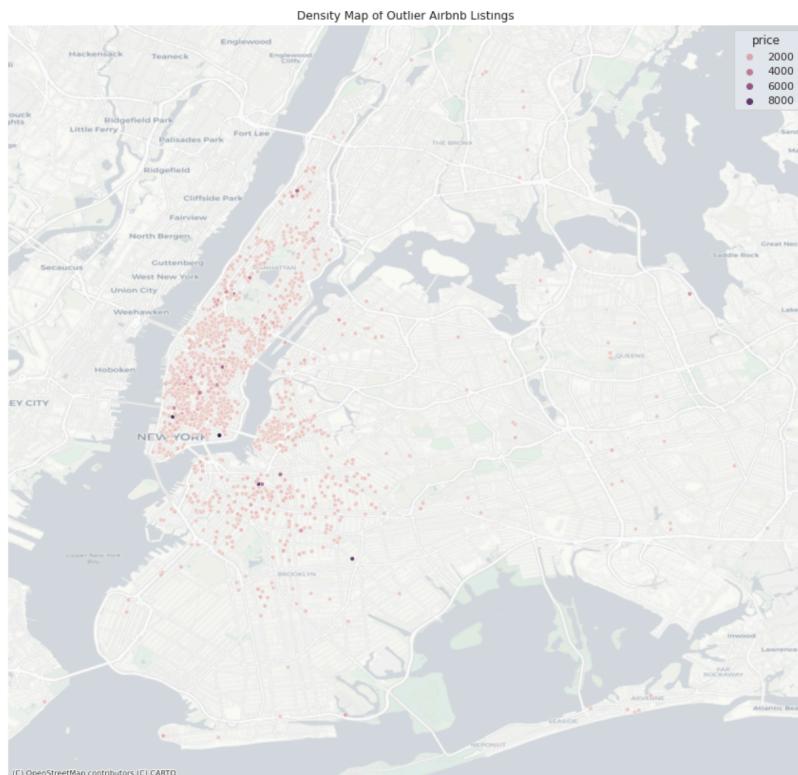
FIG. 22. (Availability based on Borough and Room

206 Types).

207 We can see a bimodal distribution with spikes at both ends. This tells us that on one
208 extreme, most NYC hosts were not leasing out their place during 2019, or maybe for a duration not
209 more than a week. On the other extreme, there were hosts that had listings available for almost the
210 entire year of 2019.

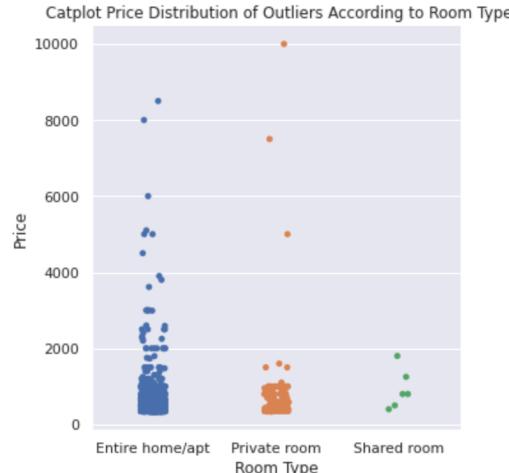
211 All boroughs had availability ranging from 0 to 365, with not much of a differentiation
212 between room types. It is interesting to note that a majority of Staten Island listings were available
213 for more days of the year compared to other boroughs.

B. Price Outliers



216 FIG. 23. (Density Map of Price Outliers of Listings).

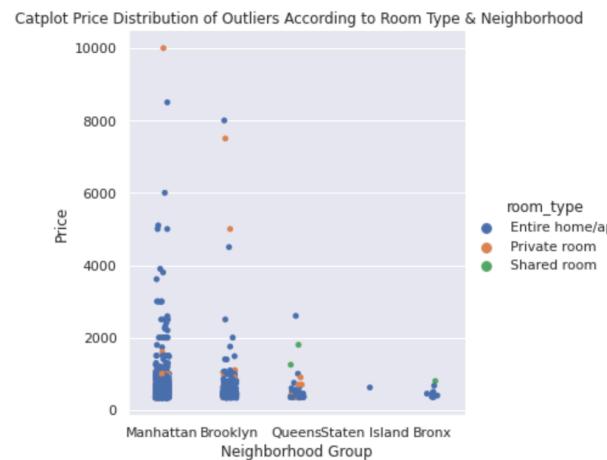
217 Most listings in our outlying price range falls around \$2000. There does not seem to be an
 218 apparent deviation of geographical correlation between expensive Airbnb listing prices compared to
 219 our tidied dataset, with a majority of the listings also concentrated in the Manhattan area. There is
 220 also only one Staten Island listing under this price range, which supports our above findings that the
 221 average Staten Island listing prices are much lower compared to the other 4 boroughs.



222

223 FIG. 24. (Catplot of Price Features

224 based on Room Type).



225 FIG. 25. (Catplot of Outlying Price Features

226 based on Room Type).

227 As seen from the plot above, most of the outlying expensive listings are either for entire
228 properties or private rooms. Interestingly, the most expensive listing is a 10000USD private room in
229 Queens. This listing is now defunct, but we can infer that it might be a long term rental situation as
230 the listing also indicates a minimum night requirement of 100 nights.

231 From this plot, it can be seen that the majority of outliers for price, seem to be under the
232 room type of an entire home/apt. The home/apt room type is also the majority for each
233 neighborhood group. We can also see that Manhattan tends to have the most number of price
234 outliers, while Brooklyn has the second highest. The logic that entire home/apt are prevalent seem
235 to make sense, as an expensive listing would be an entire home rather than a private or shared room.
236 We can also see that Staten Island tends to have the least amount of outliers that are too expensive
237 and that majority of the room type of outliers in Staten Island are entire homes/apartments.

238 **IV. NATURAL LANGUAGE PROCESSING (NLP)**

239 We will be exploring the name feature of the Airbnb dataset in depth using NLP techniques.
240 With the use of the popular NLP library Spacy we are able to clean the listing names and extract

239 relevant information from them such as the most frequently used words in listings as well as non-
240 English symbols used in the listings.

241 In our analysis of the names of the listings, we provide visualizations to show the
242 relationship of the name listings across price subgroups and popularity subgroups which we equated
243 to be the number of reviews per month.

244 Additionally, we explore the unique listing names which contain non-English symbols and
245 plot their locations to see if there is any correlation between the language used in the listing and the
246 listing location. This correlation may be present in listing names that contain Chinese characters and
247 are potentially located in Chinatown.



A word cloud visualization showing the frequency of words from the Tidydf data set. The most prominent words include 'Lower', 'East', 'space', 'Backyard', and 'block'. Other visible words include 'Brownstone', '248', '249', 'FIG.', '26.', and '(Word Cloud of Tidydf Data Set)'.

250 An initial analysis of our raw dataframe's listing names shows that the words with the most
251 frequencies are neighborhood and borough names, type of room, and general location indicators. It
252 would make sense that Brooklyn and Manhattan are some of the most frequent words because they
253 are the 2 most popular boroughs as seen in our EDA. We can also infer that hosts use listing names
254 to capitalize on different features of the property, from location indicators such as "heart" and
255 "near", number of bedrooms, and different landmarks and neighborhood names. Additionally, it is

256 interesting to note the many different variations of phrases containing "cozy", such as "cozy room",
257 "cozy bedroom", etc.

258 **A. Data Cleaning – Listing Names**

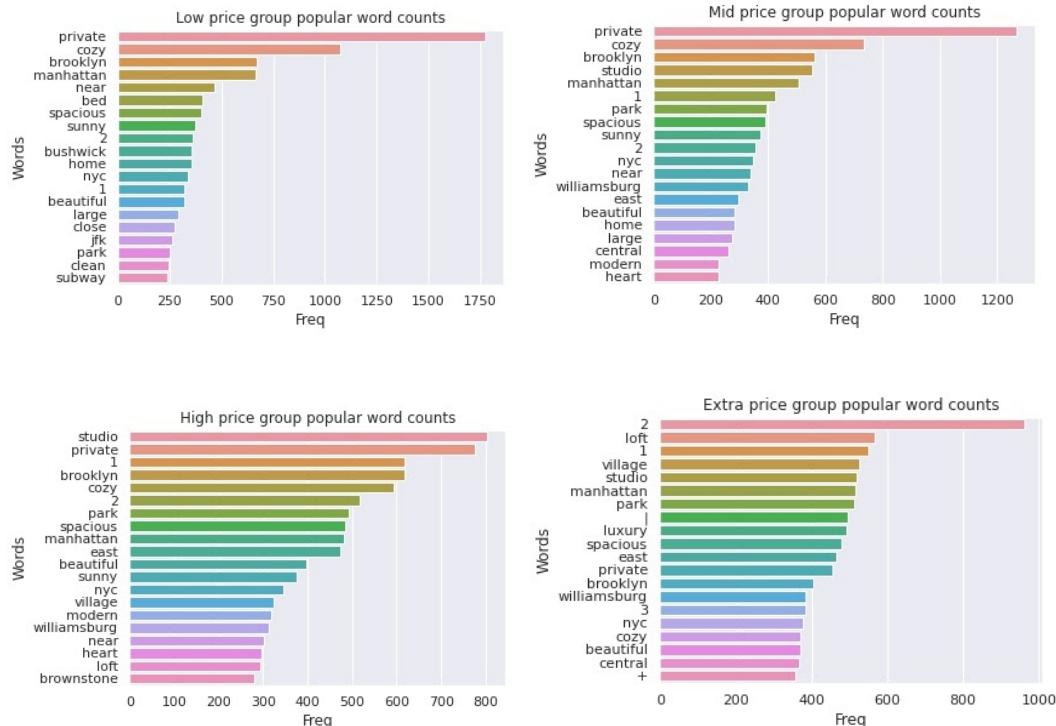
259 Apart from ensuring that we are only comparing active listings in 2019, we also observed many
260 spelling errors and large input variations which would greatly affect the quality of our NLP analysis.
261 We will thus clean and create subgroups of the listing names based on their language (English and
262 non-English) and the presence of non-alphabetical characters, particularly emojis.

263 We begin by converting our name column into a list of strings. Then using spacy's NLP package
264 we iterate through our list element of listing names and tokenize each name. We are able to identify
265 which words in each name contain punctuation or are part of spacy's stop words and we discard
266 those from each name. We also recognized that iterating through the list of names reduced our
267 computation time significantly than iterating on the name column directly. Another aspect of
268 cleaning which we kept in mind was not disregarding special characters since these characters will
269 later be used to identify non-english listings as well as listings that contain emojis. We did choose to
270 omit the words "apt", "apartment", "bedroom" from each of the names by adding them to spacy's
271 stop words since we believe that these names did not provide any relevant meaning and in our initial
272 finds were present across many subgroups of both popularity and price.

273 Here, we used NLP to get us the frequency of words within our prices category. As seen from
274 the plots below, we established price groups and got the word frequencies for each group. We
275 essentially used the summary statistics to create each group. Therefore, the low price group consists
276 of prices in the 25th percentile, medium in the 50th percentile, and high in the 75th percentile. The

277 “Extra Price” group consists of outliers in our price entries. These outliers include prices that are
278 above the 75th percentile and prices below the 25th percentile.

279 **B. Data Analysis**



280

281

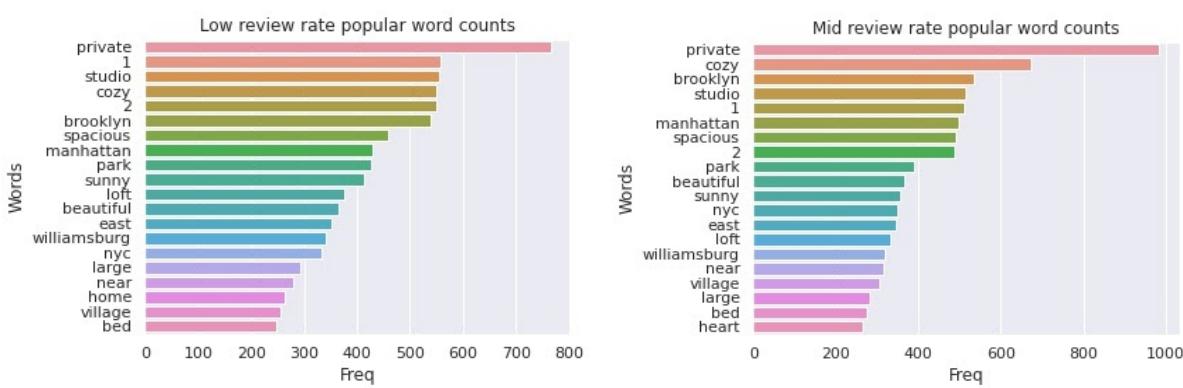
282 FIG. 27. (Count plots of Word Frequency on Type of Price).

283 From the above 4 bar charts, we can see that the term “private” is the most common, and
284 “cozy” is also the second most common term across the lower two price ranges. Referring back to
285 the word cloud, this would make sense as we concluded that the most frequent term across all
286 listings was “private room.” The prevalence of this term indicates just how important privacy is to
287 everyday Airbnb customers in New York City. This is understandable as given New York City’s high
288 population density, people would especially want a homey place to unplug from the hustle and
289 bustle.

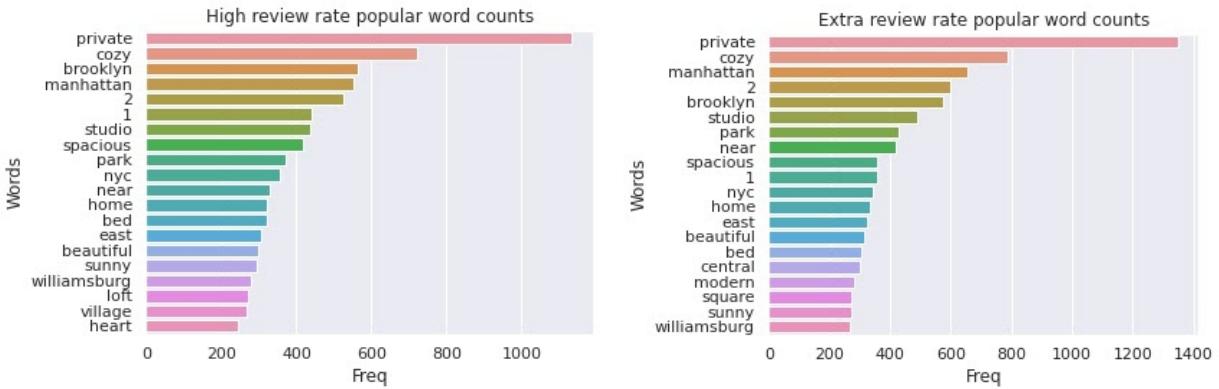
290 In comparison, many of the common terms across the higher two price brackets are of
291 property features such as “1” and “2”, most likely referring to the number of rooms or
292 “bed”[rooms], “studio”, etc. Additionally, there are terms that one would associate with more
293 expensive properties, such as “luxury”, “loft”, and “brownstone”. However, it seems like all Airbnb
294 renters want properties that are “spacious”, “beautiful”, and “sunny”.

295 There are some common geographical terms across all price ranges, such as “nyc”,
296 “manhattan”, “brooklyn”, “park”, etc. Some of these terms’ frequencies are clearly correlated with
297 price, such as “park”, most likely referring to proximity with Central Park as “central” is also a term
298 that appears in the mid and extra price range. We can also see a correlation between lower-priced
299 listings and convenience in transportation accessibility, as only the lower price bracket has the terms
300 “jfk” and “subway”. There are also specific neighborhood names that correlate to price, such as
301 “bushwick” with lower prices and “williamsburg” with higher prices.

302 In general, we can also observe that hosts with lower-priced listings prioritize appealing
303 towards everyday customers with borough names and positive descriptors, while hosts with more



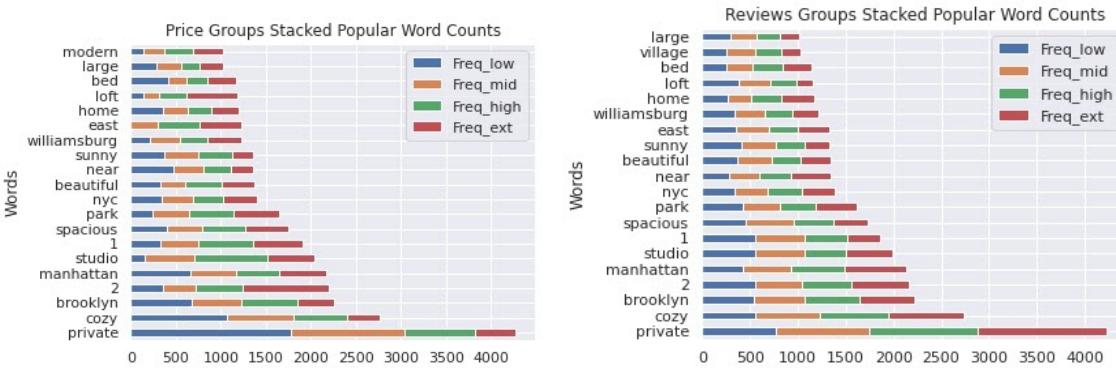
304



305

306 FIG. 28. (Count plots of Word Frequency on types of Review Rates).

307 From the above 4 bar charts, we can see that the first two most universal common terms
 308 “private” and “cozy” seem to have no effect on determining whether a renter would leave a review.
 309 Unique terms such as “modern” and “square” in listing names suggest that listings with high review
 310 rates tend to be more modern and/or mention Times Square.



311

312 FIG. 29a. (Stacked plot of Word Frequency
 313 on Price).

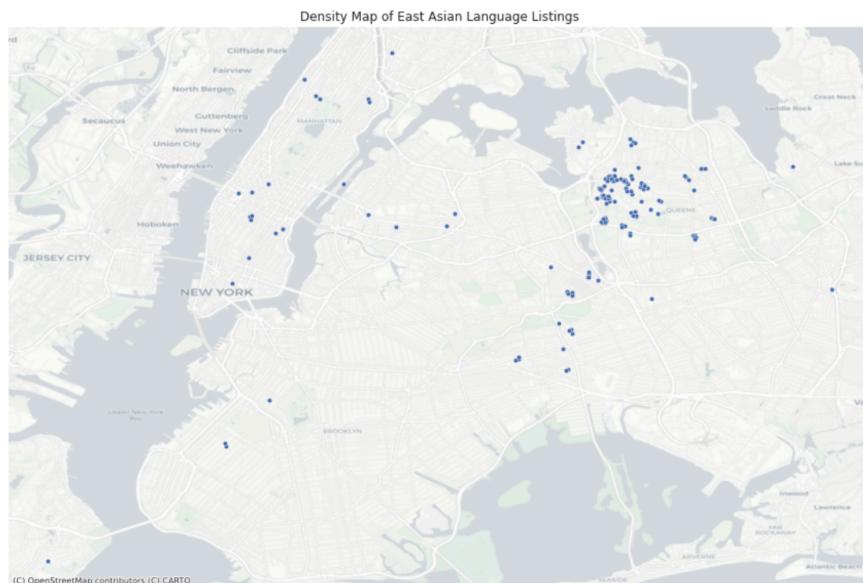
FIG. 29b. (Stacked plot of Word Frequency on
 Reviews).

314 These stacked plots above show the frequencies of the most popular words among the 4
 315 price and review groups. The stacked bar plots allow us to see differences between each group.
 316 Looking at the review groups stacked bar plots we can notice that the frequency of each word of
 317 each group looks very close to each other. While the price groups have noticeably different

318 frequencies among the same words. This is likely because there is an even distribution of review
319 scores between each of the different price categories. Looking deeper into the stacked plot for word
320 counts based on price groups, we notice that the lower-priced listings have much higher usage of the
321 words “private” and “cozy” compared to the other groups.

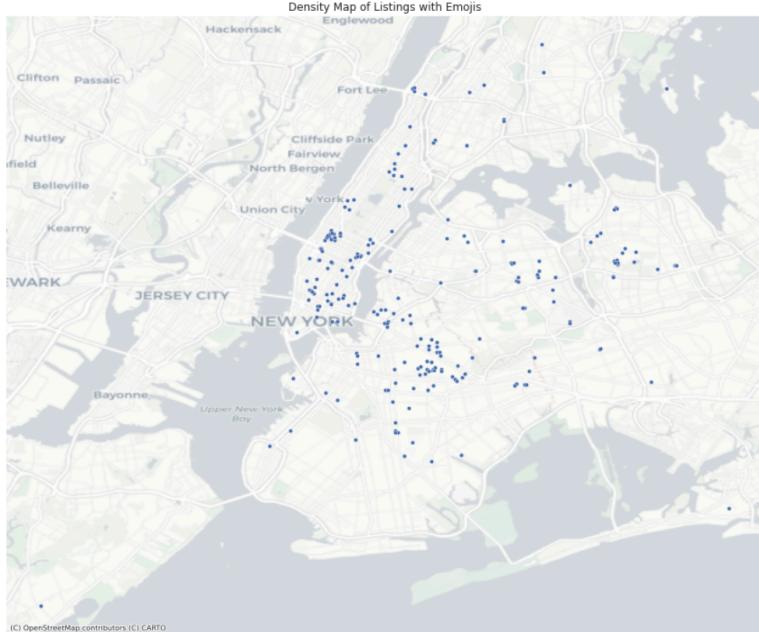
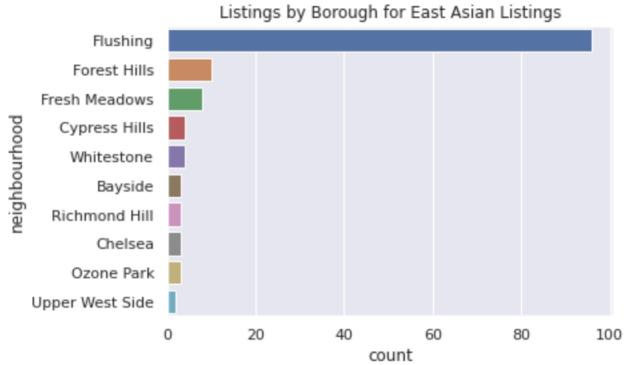
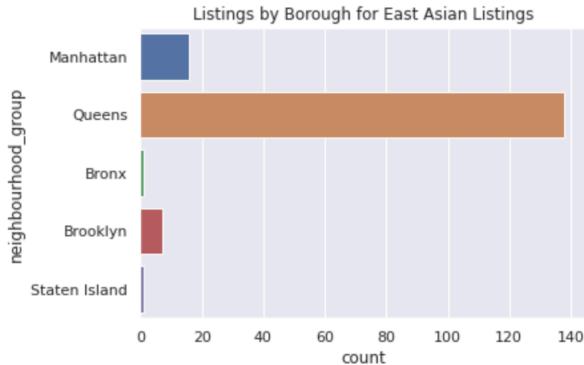
322 **C. Finding Special Characters in Data Set**

323 When looking through the list of names, we discovered certain listing names with emojis. To
324 further explore this we used simple python techniques to identify all the listing names which
325 contained characters that could not be encoded as ASCII characters. Apart from 230 names
326 containing emojis, we found names with accented letters such as é in French listings. We also found
327 160 non-English listing names in total. Among this subset was a big category of listings which
328 included East Asian characters, specifically Chinese. We reasoned that this was due to New York
329 having a prevalent East Asian community and having one of the most recognized Chinatowns in the
330 world.



331

332 FIG. 30. (Density Map of East Asian Listings).



341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

999

343 Hell's Kitchen listings show that this is purely coincidental, as there aren't any single hosts with
344 multiple emoji listings in Hell's Kitchen driving up the concentration apart from Louis (Host ID:
345 209405908) with 2 listings in Hell's Kitchen.

346 **V. CONCLUSION**

347 In our report, we have found that most New York City Airbnb listings are either in Brooklyn
348 and Manhattan, with private properties and rooms being the most popular listing type in 2019.
349 There is a very wide range of listing prices, but a majority of listings fall below \$330. There is also a
350 wide range of minimum night requirements among listings, indicating the availability of many short-
351 term and long-term rental options.

352 Through NLP, we gained further insight of how Airbnb hosts word their listing names to
353 appeal to customers. Many hosts include technical terms that state the layout of their properties,
354 geographical indicators of proximity, and positive descriptors invoking homeliness. Additionally,
355 these word choices change according to how expensive their listings are. Hosts with non-English
356 listing names, such as Chinese and other East Asian languages, attempt to appeal to corresponding
357 cultural groups based on the cultural demographic of the listings' geographic location. However,
358 there does not seem to be any immediate correlation between the use of emojis and other non-
359 ASCII characters with geography and other variables.

360 Some limitations that we encountered included the inability to accurately distinguish non-
361 English listings with English listings. Our python package only was able to differentiate between
362 listings with non-English characters. Additionally, many listing names had a lot of typos that resulted
363 in the package recognizing it as a foreign language. A suggestion for further analysis could be
364 looking into correlation of other languages with our variables apart from East Asian languages once
365 this roadblock is addressed.

366

367 APPENDIX

368 REFERENCES (BIBLIOGRAPHIC)

369 “Creating a GeoDataFrame from a DataFrame with Coordinates.” Creating a GeoDataFrame from

370 a DataFrame with Coordinates - GeoPandas 0.9.0 Documentation,

371 geopandas.org/gallery/create_geopandas_from_pandas.html.

372 Nick-Ulle. “Nick-Ulle/Teaching-Notes.” GitHub, github.com/nick-ulle/teaching-

373 notes/blob/master/sta141b/2018/discussion09.ipynb.

374 “SpaCy 101: Everything You Need to Know · SpaCy Usage Documentation.” SpaCy 101:

375 Everything You Need to Know, spacy.io/usage/spacy-101.