

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
from wordcloud import WordCloud

df = df.read\_csv("Air NYC 2019.csv")
df.head()

id name host\_id host\_name neighbourhood\_group neighbourhood latitude longitude room\_type price minimum\_nights number\_of\_reviews reviews\_per\_month availability\_365 calculated\_host\_listings\_count availability\_365

1 2539 apt hunt by the park 2787 John Manhattan Kensington 40.64749 -73.97237 Private room 149

1 2595 Skylit Midtown Castle 2845 Jennifer Manhattan Midtown 40.75362 -73.98377 Entire home/apt 225

2 3647 THE VILLAGE OF HARLEM, NYC 4632 Elisabeth Manhattan Harlem 40.80902 -73.94190 Private room 150

3 3831 Cozy Entire Floor of Brownstone 4860 Lia/Roxanne Brooklyn Clinton Hill 40.68514 -73.95976 Entire home/apt 89

4 5022 Entire Apt: Spacious Studio/Loft by central park 7192 Laura Manhattan East Harlem 40.79851 -73.94399 Entire home/apt 80

Task 1: Examine the data, there may be some anomalies in the data, and you will have to clean the data before you move forward to other tasks. Explain what you did to clean the data.

Check for zeroes in host\_id, latitude, longitude, price, minimum nights, # of reviews, availability\_365 and evaluate zeroes in host\_id, there are none (df['host\_id'] == 0).any()

Zeroes in latitude, there are none (df['latitude'] == 0).any()

Zeroes in longitude, there are none (df['longitude'] == 0).any()

Zeroes in minimum\_nights, there are none (df['minimum\_nights'] == 0).any()

Zeroes in price, there are some (df['price'] == 0).any()

(df['price'] == 0).sum()
len(df.loc[df['price'] == 0])

11

Check number of listings before
len(df)

48895

Remove zero rows because removing only 11 out of all 48895 values wouldn't change much. Values can be removed with below code
df.drop(df[df['price'] == 0].index, inplace = True)

Check number of listings after if you choose to remove 0 rows due to price
len(df)

48895

Zeroes in number\_of\_reviews, there are some, which is fine because some listings may not have been occupied (df['number\_of\_reviews'] == 0).any()

True

Zeroes in availability\_365, there are some, which is fine because some listings might be booked fully (df['availability\_365'] == 0).any()

True

Next step of clean is to look for outliers, can do so by using box plots

Boxplot for prices, there are some outliers
df.boxplot(column='price', vert=False, figsize=(20,5)).set\_title("Prices of Airbnb Listings")

Text(0.5, 1.0, 'Prices of Airbnb Listings')

Prices of Airbnb Listings

price

0 200 400 600 800 1000

Checking listings with prices over 4000. They seem to be expensive because of the location, minimum nights, and because some are luxury apartments as indicated by the name. Wouldn't make sense to remove these listings.
df.loc[df['price'] > 4000].sort\_values(by='price')

id name host\_id host\_name neighbourhood\_group neighbourhood latitude longitude room\_type

42736 33171891 30 days minimum Time square West Midtown apart... 177396569 Yanina Manhattan Hell's Kitchen 40.76043 -73.99132 Entire home/apt

28947 22296197 Chelsea Gallery for events, exhibitions fashion 3750764 Kevin Manhattan Chelsea 40.74888 -74.00481 Entire home/apt

45867 34981637 bay ridge & sunset park furnished apartment with kitchenette 263564234 Nony Brooklyn Bay Ridge 40.63087 -74.02006 Entire home/apt

43670 33796251 Beautiful private Brooklyn room with amazing view 8748976 Jeffrey Brooklyn Bedford-Stuyvesant 40.68807 -73.95426 Private room

4376 2952861 Photography Location 177497 Jessica Brooklyn Clinton Hill 40.69127 -73.95663 Entire home/apt

46614 36345358 Northside Williamsburg Turner 956324 Alex Brooklyn Williamsburg 40.71705 -73.96470 Entire home/apt

38000 30035166 4-Floor Unique Location (Greenwich Cap. -102998) 172611460 Rasmus Manhattan Harlem 40.82511 -73.94961 Entire home/apt

26739 21230553 Broadway 1 153497815 Sarah-B Brooklyn Bedford-Stuyvesant 40.68742 -73.91957 Entire home/apt

25825 20654227 Fulton 2 100089033 Sarah-2 Brooklyn Cypress Hills 40.68185 -73.88128 Entire home/apt

22353 18051877 Victorian Film Location 2675644 Alissa Staten Island Randall Manor 40.63952 -74.09730 Entire home/apt

2698 1448703 Beautiful 1 Bedroom in Nolita/Soho 213266 Jessica Manhattan Nolita 40.72193 -73.99379 Entire home/apt

4345 2919330 Near Williamsburg bridge 1211 BK Midtown 14908606 Bianca Brooklyn Bedford-Stuyvesant 40.69572 -73.95731 Private room

43009 33397385 Manhattan great location (Greenwich Cap. -102998) 16105313 Debra Manhattan Midtown 40.74482 -73.98367 Entire home/apt

3720 2243699 Penthouse Loft 1000 Loft 1483320 Omri Manhattan Little Italy 40.71895 -73.99793 Entire home/apt

3537 2110145 LWS 1BR w/bkwy&view + block from CP 26240032 Jay And Liz Manhattan Upper West Side 40.77782 -73.97848 Entire home/apt

15560 12520066 Luxury townhouse Greenwich Village 66240032 Linda Manhattan Greenwich Village 40.73046 -73.99562 Entire home/apt

29684 22780103 Park Avenue Mansion by (Hidden by Airbnb) 156158778 Sally Manhattan Upper East Side 40.78517 -73.95270 Entire home/apt

3774 2271504 Brooklyn Duplex 11598359 Jonathan Manhattan Clinton Hill 40.68766 -73.96439 Entire home/apt

37194 29547314 Apartment New York (White's Kitchens 35303743 Patricia Manhattan Upper West Side 40.76835 -73.98367 Private room

48043 36056808 Luxury Tribeca Apartment at 271248669 Jenny Manhattan Tribeca 40.71206 -74.00999 Entire home/apt

44034 33998396 3000 sq ft daylight photo studio 3750764 Kevin Manhattan Chelsea 40.75060 -74.00388 Entire home/apt

42523 33007610 70 Luxury MotorYacht on the Hudson 7407743 Jack Manhattan Battery Park City 40.71162 -74.01693 Entire home/apt

45666 34895693 Gem of east Flatbush 262534951 Sandra Brooklyn East Flatbush 40.65724 -73.92450 Private room

29862 22797276 East 72nd Townhouse by (Hidden by Airbnb) 156158778 Sally Manhattan Upper East Side 40.76824 -73.95989 Entire home/apt

4377 2953058 Film Location 1774797 Jessica Brooklyn Clinton Hill 40.69137 -73.96723 Entire home/apt

30268 23377410 Beautiful/Spacious 1br luxury flat-TriBeCa/Soho 18128455 Rum Manhattan Tribeca 40.72197 -74.00633 Entire home/apt

4043 31340283 2br + The Heart of NYC Manhattan Lower East... 4382127 Matt Manhattan Lower East Side 40.71980 -73.98566 Entire home/apt

6530 4737930 Spanish Harlem Apt 1235070 Olson Manhattan East Harlem 40.79264 -73.93898 Entire home/apt

12342 9528920 Quiet, Clean, Lit LES & Chinatown 3906464 Amy Manhattan Lower East Side 40.71355 -73.98507 Private room

29238 22436899 1-BR Lincoln Center 72390391 Jelena Manhattan Upper West Side 40.77213 -73.98665 Entire home/apt

9151 7003697 Furnished room in Astoria apartment 20582832 Kathrine Queens Astoria 40.76810 -73.91651 Private room

17892 13894339 Luxury 1 bedroom apt - stunning Manhattan view 5143901 Erin Brooklyn Greenpoint 40.73260 -73.95739 Entire home/apt

Task 2: Examine how the prices of the Airbnb changes with the change in the neighborhood.
a. Find Top 5 and bottom 5 neighborhood based on the price of the Airbnb in that neighborhood (select only neighborhoods with more than 5 listings). (10 Points)
b. Analysis, the price variation between different neighborhood group, and plot these trends. (5 Points)

Top 5 neighborhood based on prices with neighborhoods with more than 5 listings are:
Williamsburg, Bedford-Stuyvesant, Harlem, Bushwick and Upper West Side
bot5df = df['neighbourhood'].value\_counts()
top5df = top5df[top5df > 5]
top5df

Williamsburg 3920
Bedford-Stuyvesant 3114
Harlem 2458
Bushwick 2465
Upper West Side 1971
..
Bull's Head 6
Midland Beach 6
Grant City 6
Mount Eden 6
Bay Terrace 6
Name: neighbourhood, dtype: int64

Williamsburg 3920
Bedford-Stuyvesant 3114
Harlem 2458
Bushwick 2465
Upper West Side 1971
..
Bull's Head 6
Midland Beach 6
Grant City 6
Mount Eden 6
Bay Terrace 6
Name: neighbourhood, dtype: int64

Part B, check price variation between different groups
ngdf = df.groupby('neighbourhood\_group')['price'].mean().to\_frame()
ngdf.plot.bar(y='price', rot=0, title='Median of Neighborhood Groups in Airbnb Homes')

Mean of Neighborhood Groups in Airbnb Homes

price

0 25 50 75 100 125 150 175 200

Bronx Brooklyn Manhattan Queens Staten Island

neighbourhood\_group

ngdf = df.groupby('neighbourhood\_group')['price'].median().to\_frame()
ngdf.plot.bar(y='price', rot=0, title='Median of Neighborhood Groups in Airbnb Homes')

Median of Neighborhood Groups in Airbnb Homes

price

0 20 40 60 80 100 120 140

Bronx Brooklyn Manhattan Queens Staten Island

neighbourhood\_group

ngdf = df.groupby('neighbourhood\_group')['price'].std().to\_frame()
ngdf.plot.bar(y='price', rot=0, title='Standard Deviation of Neighborhood Groups in Airbnb Homes')

Standard Deviation of Neighborhood Groups in Airbnb Homes

price

0 50 100 150 200 250 300

Bronx Brooklyn Manhattan Queens Staten Island

neighbourhood\_group

Task 3: Select a set of the most interesting features. Do a pairwise Pearson correlation analysis on all pairs of these variables. Show the result with a heat map and find out most positive and negative correlations. (5 points)

df.corr()

id host\_id latitude longitude price minimum\_nights number\_of\_reviews reviews\_per\_month availability\_365 calculated\_host\_listings\_count

id 1.000000 0.588290 -0.003125 0.090908 0.010619 -0.013224 -0.319760 -0.017600 -0.000116 -0.000000
host\_id 0.588290 1.000000 0.020224 0.127055 0.015309 -0.017364 -0.140139 -0.015389 -0.015389
latitude -0.003125 0.020224 1.000000 0.084788 0.033939 0.024869 -0.030608 -0.021702 1.000000
longitude 0.090908 0.127055 0.084788 1.000000 -0.150019 -0.062747 0.005904 -0.047954 -0.030608
price 0.010619 0.015309 0.033939 -0.150019 1.000000 0.042799 -0.047954 -0.080116 -0.000000
minimum\_nights -0.013224 -0.017364 -0.024869 -0.062747 0.042799 1.000000 -0.080116 -0.000000
reviews\_per\_month -0.319760 -0.140139 -0.015389 0.059094 -0.047954 -0.080116 1.000000
calculated\_host\_listings\_count 0.000000 0.296417 -0.010142 0.145948 -0.030608 -0.121702 0.549868 -0.072376
availability\_365 0.085468 0.203492 -0.01093 0.062731 0.081629 0.144303 0.172028

Interesting Features
# Most positive correlations:
# Reviews per month and number of reviews (0.55). This makes sense because the number of reviews per month is the total number of reviews divided by the number of months a property was listed for.
# Correlation Coefficient between Neighborhood Group and Room Type = 0.66
# Minimum nights and reviews per month
corrdf = df.corr()
corrdf = corrdf.drop\_duplicates(subset=['id', 'host\_id', 'latitude', 'longitude'])
dataplot = sb.heatmap(corrdf.corr(), cmap='YlGnBu', annot=True).set\_title('Correlations between Interesting Listing Features')
plt.set(rc={'figure.figsize':(8,8)})
plt.show()

Correlations between Interesting Listing Features

latitude longitude price minimum\_nights number\_of\_reviews reviews\_per\_month calculated\_host\_listings\_count availability\_365

latitude 1.0000 0.095 0.034 0.025 0.019 0.01 0.02 -0.011
longitude -0.085 1.000 -0.15 -0.063 -0.059 0.15 0.11 0.083
price -0.034 -0.15 1.000 0.043 -0.048 0.031 0.057 0.082
minimum\_nights -0.025 -0.063 0.043 1.000 0.08 0.12 0.13 0.014
number\_of\_reviews -0.015 0.095 0.048 0.08 1.000 0.55 0.072 0.017
reviews\_per\_month -0.01 0.15 -0.031 0.12 0.55 1.000 0.004 0.19
calculated\_host\_listings\_count -0.011 0.057 0.13 -0.072 0.004 1.000 1.000
availability\_365 -0.011 0.083 0.082 0.14 0.17 0.19 0.23 1.000

Task 4: The latitude and longitude of all the Airbnb listings are represented in the dataset.
a. Plot a scatter plot based on these coordinates, where the points represent the location of an Airbnb, and the points are color coded based on the neighborhood group feature. (5 Points)
b. Now again, plot based on these coordinates, where the points represent the location of an Airbnb, and the points are color coded based on the price of the particular Airbnb, where price of the listing is less than 1000. Looking at the graph can you tell which neighborhood group is the most expensive. (10 Points)

Part A
sb.scatterplot(data=df, x='latitude', y='longitude', hue='neighbourhood\_group').set\_title('Airbnb Listings Based on Location')
plt.xlim(-74.3, -73.7)
plt.ylim(40.4, 41)

Airbnb Listings Based on Neighborhood Groups

neighbourhood\_group

Brooklyn
Manhattan
Queens
Staten Island
Bronx

Part B
After looking at the graph, it is clear that the most expensive location seems to be Manhattan as the area where the most expensive listings are correspond with where Manhattan is, as seen in the previous graph.
pricedf = df.copy()
pricedf = pricedf.loc[pricedf['price'] < 1000].sort\_values(by='price')
sb.scatterplot(data=pricedf, x='latitude', y='longitude', hue='price').set\_title('Airbnb Listings Based on Price')
plt.xlim(-74.3, -73.7)
plt.ylim(40.4, 41)

Airbnb Listings Based on Price

price

0
200
400
600
800

Task 5:
Word clouds are useful tool to explore the text data. Extract the words from the name of the Airbnb and generate a word cloud.

list = df['name'].to\_frame().stack().str.split('\n+').explode().tolist()
list2 = list
for i in list:
list2.append(i.replace(' ', ''))
print(list2[0])
wordcloud = WordCloud.generate\_from\_list(list2)
wordcloud = WordCloud(width = 1000, height = 1000).generate\_from\_list(list2)
plt.figure()
plt.imshow(wordcloud)
plt.axis('off')
plt.show()

Task 6: Find out which areas has the busiest (hosts with high number of listings), price? Are there any reasons, why these areas are the busiest, considers factors such as availability, host, review, etc.? Bolster your reasoning with different plots and correlations. (10 Points)

Part A
The first scatter plot is made to see the locations where hosts have a lot of listings. The locations can be eyeballed with the graph scatter plots that were made in Part 4 in order to see the neighborhood groups based on the latitude and longitude of the locations.
busydf = df.copy()
busydf = busydf.loc[busydf['calculated\_host\_listings\_count'] > 30].sort\_values(by='calculated\_host\_listings\_count')
sb.scatterplot(data=busydf, x='latitude', y='longitude', hue='calculated\_host\_listings\_count').set\_title('Airbnb Listings Based on Location')
plt.xlim(-74.3, -73.7)
plt.ylim(40.4, 41)

Airbnb Listings Based on Neighborhood Groups

calculated\_host\_listings\_count

50
100
150
200
250
300

Part B
In order to see the correlation between the busiest locations, the listings with owners that had few listings were removed as a lot of those listings wouldn't be super representative of the busiest areas.
corrdf = df.corr()
corrdf = corrdf.loc[corrdf['calculated\_host\_listings\_count'] > 30].sort\_values(by='price')
corrdf = corrdf.drop\_duplicates(subset=['host\_id'], keep='first')
corrdf = corrdf.drop(columns=['id', 'host\_id', 'latitude', 'longitude'])
corrdf['neighbourhood\_group'] = corrdf['neighbourhood\_group'].map({'Staten Island':0, 'Bronx':1, 'Queens':2, 'Brooklyn':3, 'Manhattan':4})
corrdf['room\_type'] = corrdf['room\_type'].map({'Shared room':0, 'Private room':1, 'Entire home/apt':2})
dataplot = sb.heatmap(corrdf.corr(), cmap='YlGnBu', annot=True).set\_title('Correlations between Interesting Listing Features')
plt.set(rc={'figure.figsize':(8,8)})
plt.show()

Correlations between Interesting Listing Features for Busiest Locations

neighbourhood\_group room\_type price minimum\_nights number\_of\_reviews reviews\_per\_month calculated\_host\_listings\_count availability\_365

neighbourhood\_group 1.000 0.66 0.4 0.00079 -0.0057 0.089 0.24 0.3
room\_type 0.66 1.000 0.48 0.18 -0.13 -0.057 0.29 0.16
price 0.4 0.48 1.000 -0.055 -0.11 0.25 0.25 0.025
minimum\_nights 0.00079 0.18 -0.055 1.000 -0.42 -0.75 -0.27 -0.038
number\_of\_reviews -0.0057 -0.13 -0.11 -0.42 1.000 0.5 0.016 0.044
reviews\_per\_month 0.089 -0.057 0.25 0.75 0.5 1.000 0.7 0.06
calculated\_host\_listings\_count 0.24 0.29 0.25 -0.27 -0.016 0.7 1.000 0.035
availability\_365 0.3 0.16 0.025 -0.038 0.044 0.06 0.035 1.000

corrdf = df.corr()
corrdf = corrdf.loc[corrdf['calculated\_host\_listings\_count'] > 30].sort\_values(by='price')
corrdf = corrdf.drop(columns=['id', 'host\_id', 'latitude', 'longitude'])
corrdf['neighbourhood\_group'] = corrdf['neighbourhood\_group'].map({'Staten Island':0, 'Bronx':1, 'Queens':2, 'Brooklyn':3, 'Manhattan':4})
corrdf['room\_type'] = corrdf['room\_type'].map({'Shared room':0, 'Private room':1, 'Entire home/apt':2})
dataplot = sb.heatmap(corrdf.corr(), cmap='YlGnBu', annot=True).set\_title('Number of Listings for Busiest Hosts Based on Location')
plt.set(rc={'figure.figsize':(8,8)})
plt.show()

Number of Listings for Busiest Hosts Based on Location

Queens Bronx Manhattan Brooklyn

Task 7: Create two plots (at least one unique plot not used above) of your own using the dataset that you think reveals something very interesting. Explain what it is, and anything else you learned. (10 Points)

Part 6 Analysis:
After looking at the scatter plot and the previous scatter plot in Part 4, it seems that most of the busiest listings seem to be in Manhattan. This fact is later confirmed by the bar chart for the Number of Listings for Busiest Hosts Based on Location. When looking at the correlation analysis for the busiest hosts, there are some numbers that stand out:
1. Correlation Coefficient between Neighborhood Group and Room Type = 0.66
2. Calculated Host Listings Count and Reviews Per Month = 0.7
3. Neighborhood Group and Price = 0.4
4. Room Type and Price = 0.48
The coefficient between Neighborhood Group and Room Type is interesting because when thinking about it, whether a listing is shared, private or the entire home/apartment does relate to the price. In the 5 boroughs, Manhattan can have both entire apartments or private or shared rooms, while in places like Queens, you're probably getting the entire home/apartment. Another value that was very intriguing was the coefficient between Calculated Host Listings Count and Reviews Per Month because it is telling us that with the busiest hosts, the more listings they have, the more reviews they are getting. With the last two values, it makes sense that the location and room type affects how costly a listing is.

First Unique Plot: This is created as a pie plot, that is used as a proportional representation of data in a column. Pie plots take in numeric data and the neighborhood column is of a different type, so the number of occurrences were counted in order to quantify this data. Since the number of listings are well over over 40000 listings, I wanted to find the most popular neighborhoods and see which places have the most listings within those areas. I selected locations that have more than 1000 listings and saw that Williamsburg and Bedford-Stuyvesant, which are both in Brooklyn, were the most popular neighborhoods in the data set.
top5df = df['neighbourhood'].value\_counts()
top5df = top5df[top5df > 1000]
plt.pie(top5df.autotopct=41.05%).set\_title('Most Popular Neighborhoods Based on Airbnb Listings')
Text(0.5, 1.0, 'Most Popular Neighborhoods Based on Airbnb Listings')

Most Popular Neighborhoods Based on Airbnb Listings

Bedford-Stuyvesant 14%
Williamsburg 15%
Chelsea 4%
Greenpoint 4%
East Harlem 4%
Midtown 6%
Crown Heights 8%
Upper East Side 7%
East Village 7%
Hell's Kitchen 7%
Upper West Side 7%
Bushwick 9%
Harlem 10%

Another element I wanted to reveal about the data set was the types of rooms that were in the Airbnb listings, based on the the location of the listing. I used a scatter plot to do so After creating the chart (and using a previous chart from Part 4 for the location based on the latitude and longitude), I was able to see that in Bronx and in Brooklyn there were more private rooms, in Queens, there was a mix between private rooms and entire homes/apartments, and in Manhattan there were mostly listings that included the entire home/apartment.

corrdf = df.corr()
corrdf = corrdf.drop(columns=['id', 'host\_id', 'latitude', 'longitude'])
corrdf['neighbourhood\_group'] = corrdf['neighbourhood\_group'].map({'Staten Island':0, 'Bronx':1, 'Queens':2, 'Brooklyn':3, 'Manhattan':4})
corrdf['room\_type'] = corrdf['room\_type'].map({'Shared room':0, 'Private room':1, 'Entire home/apt':2})
dataplot = sb.heatmap(corrdf.corr(), cmap='YlGnBu', annot=True).set\_title('Airbnb Listings Based on Room Type')
plt.xlim(-74.3, -73.7)
plt.ylim(40.4, 41)

Airbnb Listings Based on Room Type

room\_type

Private room
Entire home/apt
Shared room