

Capstone Project I

Milestone Report

Problem Statement

Predicting the likelihood of each Airbnb host to cancel on their guests. Along with the price and ratings shown in their listings, an Airbnb host's likelihood of cancelling should also be readily available for prospective guests to see.

As part of the "share economy", Airbnb hosts become the key determinants of the guest's experience with the company. This experience can determine whether or not the guest will continue to use the app in the future.

The proposed added feature will give the guests more information about the host. This feature will help the guests decide whether or not they would like to take the risk, knowing the host's likelihood to cancel (the effectiveness of the added metric, of course, will depend on the guest's risk preference).

This can also become an incentive for the host to cancel on guests less (right now, they only get a \$50/\$100 penalty when they cancel).

Data Wrangling Report

In this project, I focus on the Airbnb market in London, the second largest Airbnb city outside the US.

Most of the datasets used in this project were all found on [insideairbnb.com](#). This website "is an independent, non-commercial set of tools and data that allows you to explore how Airbnb is really being used in cities around the world."

The following are the csv files found on the Inside Airbnb website:

- listings.csv
 - information and metrics for listings in London
- reviews.csv
 - reviews left by the guests on the hosts' listings
- neighbourhoods.csv
 - Lists all neighborhoods and neighborhood groups in London

I also want to see the effect that the Airbnb listing's neighborhood's crime rate has on our study so I also make use of the London crime report posted in the 'London Datastore' website. This csv contains all the crime rates in all London boroughs. The data dates back from 1999.

listings.csv

columns:

id	room_type
name	price
host_id	minimum_nights
host_name	number_of_reviews
neighbourhood_group	last_review
neighbourhood	reviews_per_month
latitude	calculated_host_listings_count
longitude	availability_365

The listings data contains 80,767 rows and 16 columns. Each row represents a unique listing in London.

This is one of our main datasets as it gives us a lot of information about the listings. When we take the minimum of `number_of_reviews`, we get 0. This is relevant in this study because we get the information on the cancellations through the automated postings listed on the reviews section. We cannot simply set the number of cancellations in a listing as 0 because of 0 reviews; this could

also mean that the listing is inactive. We determine inactivity by checking those airbnbs' `availability_365` and see if any of them have 365 days available. Once the inactive listings are determined, they are dropped from the listings DataFrame.

reviews.csv **columns:**

listing_id	As previously mentioned, we track cancellations made by the listings by looking at the
id	<code>comments</code> . When a host cancels on a guest, an automated review is posted on their profile,
date	indicating that they canceled on the guest. This review cannot be deleted. This reviews data
reviewer_id	contains 1,249,466 rows and 6 columns--where each row represents a comment written on a
reviewer_name	specific listing's page. When we apply <code>.info()</code> to the DataFrame, we see that there are 407 null
comments	values (i.e. empty reviews) on the "comments" column. Since this is only .03% of the reviews,
	we can drop these rows using <code>.dropna()</code> .

neighbourhoods.csv

The neighborhood data is simple and just gives us the 33 neighborhoods in London. There are no null values or outliers.

crimerates.csv

Code	In this dataset, we are only interested in:
Borough	<ul style="list-style-type: none">Overall crime rate of each neighborhood (categorized as "All recorded offences" under the <code>Offences</code> column), so we drop the other offences
Year	<ul style="list-style-type: none">Neighborhoods included in the <code>neighbourhood</code> DataFrame. Places like "Inner London", "England and Wales", "Met Police Area", "Outer London", and "Heathrow" are
Offences	included in the <code>Borough</code> column despite not being an official borough (i.e. neighborhood)
Rate	in London and thus not included in our <code>neighbourhood</code> DataFrame. We simply ignore this
Number_of_offences	and drop it from our data.

- The crime rates from 2011 to present since the earliest data we have on our Airbnb data is from 2011. We do this by dropping all the crime rates from the years prior. We use:

```
all_crimes[(all_crimes['Year'])>=2011]
```

Creating the Main DataFrame

Using all these datasets, we create a main DataFrame where each row represents a unique Airbnb listing. Besides merging the DataFrames, new columns are created to provide better understanding of our data.

In the `comments` tab of the reviews data, there were two versions of the automated postings: "The host canceled my reservation n days before arrival" and "The reservation was canceled n days before arrival. This is an automated posting." Upon further review, we see that the automated posting template changed sometime in 2012. We create a new DataFrame only containing listings with automated postings in their comments, and apply `.value_counts()` on the `comments` tab. This will give us information on how many times reservations for each listing were cancelled before. We add this in our main DataFrame as a new tab, `num_cancellations`.

We also add a `rate` tab to the DataFrame. The rate represents the overall crime rate of the neighborhood that the listing belongs in.

After all the data wrangling, we come up with a main DataFrame with the following columns:

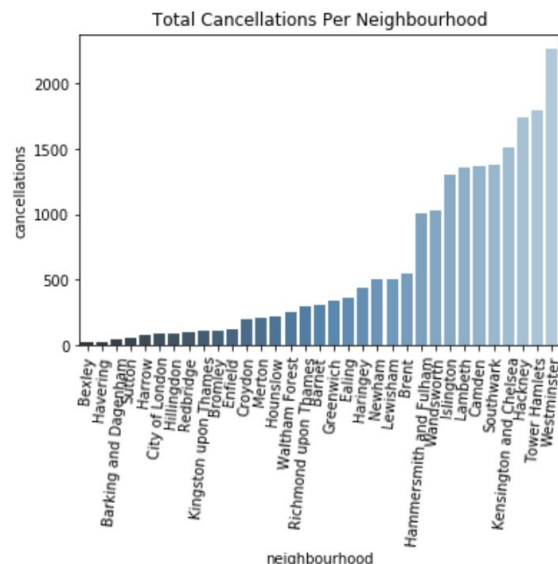
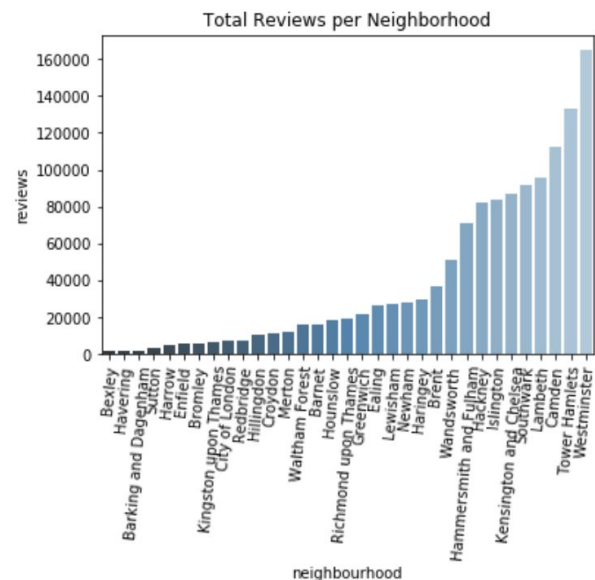
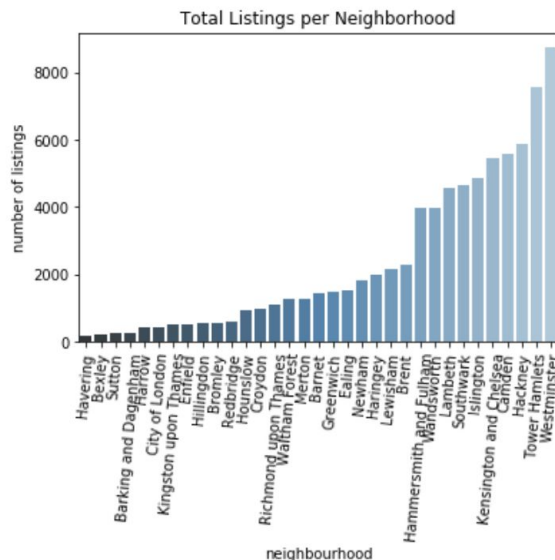
neighbourhood	days booked
room_type	num_cancellations
price	Year
minimum_nights	Rate
number_of_reviews	Number_of_offences
calculated_host_listings_count	

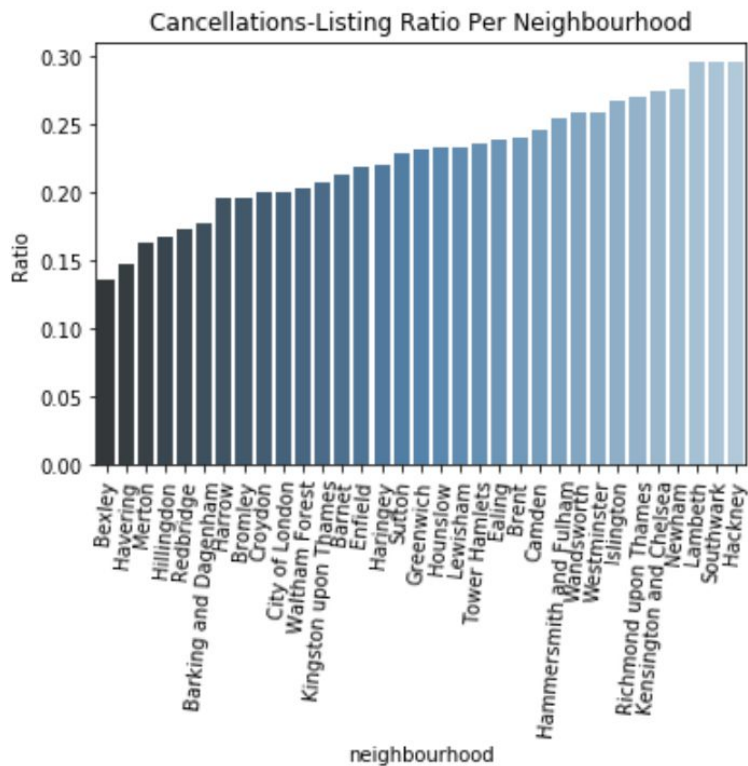
Exploratory Data Analysis

Next step is to ask questions and look into possible trends between the features of the Airbnb listings.

The following barplots were created to see the distribution of listings, reviews and cancellations across the London neighborhoods.

Which neighborhoods have the most Airbnb listings? Which neighborhoods have the most Airbnb reviews? Which neighborhoods have the most Airbnb cancellations?

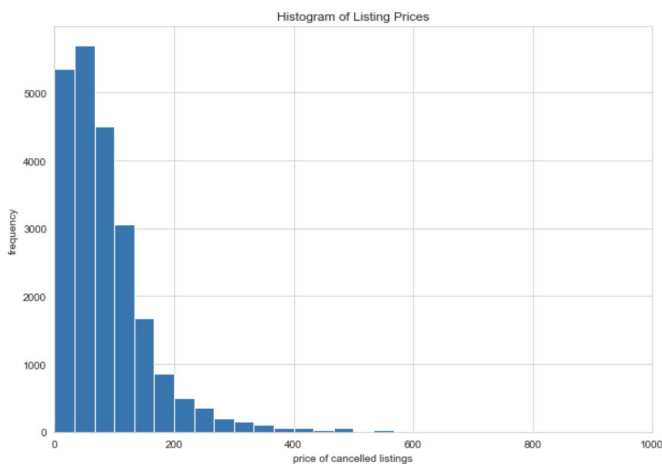




Westminster and Tower Hamlets are the top two neighborhoods with the highest frequency in all features--could this mean that listings in their neighborhoods have the highest likelihood of cancelling or is it just because they also have the most number of listings? This question is addressed by creating a barplot that shows the ratio of the number of cancellations to the number of listings per neighborhood.

This barplot shows us that 1) the distribution is no longer as spread out and 2) Westminster does not have the highest ratio.

Does price correlate with these cancellations?

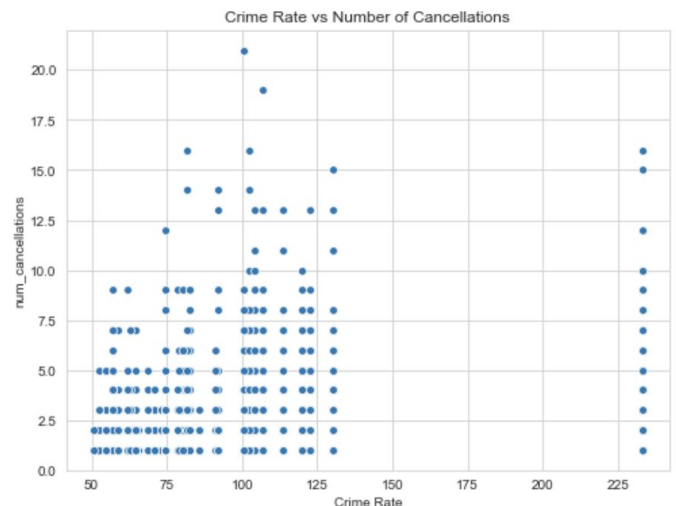


A histogram of the price of each cancelled listing reservation is created to study this relationship.

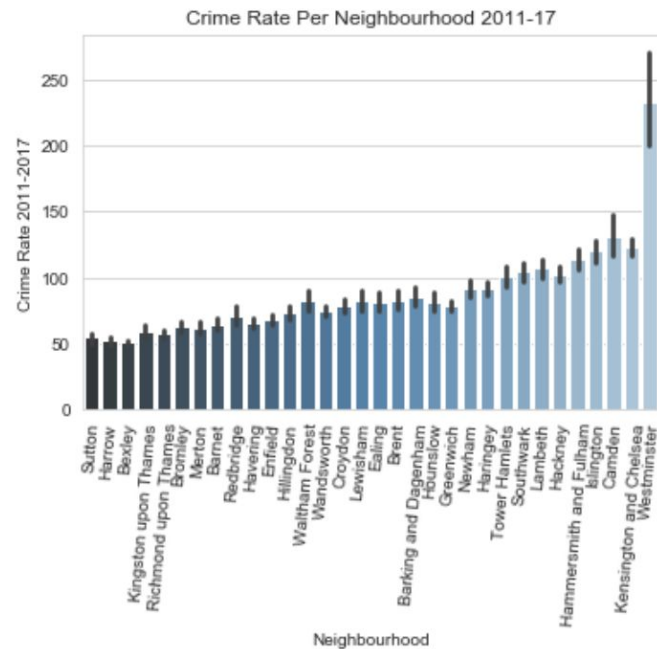
The histogram shows a negative correlation between the two variables.

Does the neighborhood crime rate correlate with the cancellations in the area?

The relationship between the two seems random. We can confirm if the correlation is indeed insignificant by using statistical techniques later.



The heatmap shows crime has gone down for all of London since 1990. Westminster had the biggest change, although, it still remains to have the highest crime rate.



Since the crime rates of all neighborhoods have decreased since 1990, we can just focus on the crime rates from 2011 (earliest Airbnb listing in our data) to 2017 (latest crime rate data provided by the csv file).

Inferential Statistics

After studying trends between the Airbnb listings, we apply inferential statistics techniques to see if any of these relationships are statistically significant. We mainly focus on how these features are related to the number of cancellations of the listings.

For most of the variables, we make use of permutation tests to see if their respective correlations with the number of cancellations are statistically significant.

Cancellations vs. Price

Null Hypothesis: There is no significant correlation between the number of cancellations and the price of the listing.

Alternative Hypothesis: There is a significant correlation between the two.

First, we calculate the actual correlation between the number of cancellations and listing prices. The observed correlation is -0.02526.

We test to see if this observed correlation is statistically significant by performing a permutation test. We use `np.random.permutation` to reorder the listing prices and get their correlation with the number of cancellations each time.

In this case, our p-value is the ratio of the amount of times the absolute value of our correlation replicates was greater or equal to our observed correlation.

```
observed correlation between price and number of cancellations: -0.02526
p-value: 0.0000000000
We reject the null. There is a significant correlation between price and
number of cancellations
```

The p-value we get is really small. We then conclude that the correlation between price and the number of cancellations is statistically significant at the 1% level.

Cancellations vs Popularity (Number of Reviews)

Null Hypothesis: There is no significant correlation between the number of cancellations and the number of reviews a listing gets.

Alternative Hypothesis: There is a significant correlation between the two.

We perform the same test and conclude that their observed correlation of 0.21254 is statistically significant at the 1% level.

```
observed correlation between number of reviews and number of cancellations: 0.21254
p-value: 0.0000000000
We reject the null. There is a significant correlation between the number of cancellations and number of reviews
```

Cancellations vs Demand (Days Booked)

Null Hypothesis: There is no significant correlation between the number of cancellations and the number of days the listing is booked.

Alternative Hypothesis: There is a significant correlation between the two.

We perform the same test and conclude that their observed correlation of -0.01100 is statistically significant at the 1% level.

```
observed correlation between number of days booked and number of cancellations: -0.01100
p-value: 0.0018000000
We reject the null. There is a significant correlation between the number of cancellations and the number of days the listing is booked.
```

Cancellations vs Minimum Nights

Null Hypothesis: There is no significant correlation between the number of cancellations and the minimum nights required by a listing.

Alternative Hypothesis: There is a significant correlation between the two.

We perform the same test and conclude that their observed correlation of -0.00473 is **not** statistically significant at the 1% level.

```
observed correlation between minimum nights and number of cancellations: -0.00473
p-value: 0.1630000000
We fail to reject the null. There is no significant correlation between the number of cancellations and the minimum nights required by a listing.
```

Cancellations vs Crime

Null Hypothesis: There is no significant correlation between the number of cancellations and the crime rate of the neighbourhood that the listing belongs in.

Alternative Hypothesis: There is a significant correlation between the two.

We perform the same test and conclude that their observed correlation of -0.01161 is statistically significant at the 1% level.

```
observed correlation between crime rate and number of cancellations: 0.01161
p-value: 0.0009000000
We reject the null. There is a significant correlation between the number of cancellations and the crime rate of the neighbourhood that the listing belongs in.
```

Besides these features, we also want to check if the number of cancellations vary by categories such as neighborhood and room type. Since there are more than two in each category, we cannot perform ttests or permutation tests. Instead, we can use an F-test. In an F-test, we are able to compare the means of various groups and determine if they are equal by looking at their variances. With this technique, we are able to check variations between and within the groups.

Cancellation by Neighborhood?

Null Hypothesis: The average number of cancellations of each neighborhood are not statistically different from each other.

Alternative Hypothesis: The average number of cancellations of each neighborhood are statistically different from each other.

```
f-stat: 4.65329
p-value: 0.0000000000
The average number of cancellations of each neighbourhood are statistically different from each other.
```

The f-stat of 4.65 suggests that between-neighbourhood variance is 4.65 times the within-neighborhood variance. The small p-value suggests that the difference is statistically significant.

Cancellation by Room Type?

Null Hypothesis: The average number of cancellations of each room type are not statistically different from each other.

Alternative Hypothesis: The average number of cancellations of each room type are statistically different from each other.

```
f-stat: 11.68846
p-value: 0.0000084047
The average number of cancellations of each room type are statistically different from each other.
```


The p-value suggests that the differences in the number of cancellations among the room type are also statistically significant.

Summary Findings

After studying trends through data visualization and applying inferential statistics techniques, we conclude that a listing's number of cancellations has significant correlation with the following variables:

- Price
- Reviews
- Availability
- Crime rate

We also conclude that the number of cancellations vary significantly across neighborhoods and room type.