

NYC Airbnb House Price vs Crime

Team member name 1: Keye Chen (keyechen); section (001)

Team member name name 2: Fengyu Zhang (fyzhang); section (001)

Team member name name 3: Boyan Wu (wuboyan); section (001)

1. Motivation

1.1 Abstract

When people choose their accommodations, they weigh numerous factors, with cost and safety topping the list. In this project, we are delving into the dynamics of Airbnb pricing and crime across New York City's vibrant neighborhoods

1.1.1 Demands

- Our inquiry is driven by the need to understand how Airbnb prices are influenced by crime rates, Concurrently, the project aims to equip renters with a knowledge base to make informed accommodation choices, and city planners and law enforcement to understand the relationship between tourism accommodations and crime, effectively bridging the gap between cost and security.

1.1.2 Questions

- The central question guiding this research is: what variations in Airbnb house prices can be partly attributed to the perceived safety or crime statistics of different neighborhoods? What are the driving factors behind our analysis and how can these insights be utilized to optimize the findings and performance around geographic regions in NYC?

1.1.3 Goal

- This study aims to answer questions like the nature of the correlation between Airbnb house prices and crime rates, identify neighborhoods with significant discrepancies, and predict trends in these metrics. The goal is to provide actionable insights and recommendations for various stakeholders.

2. Data Sources

2.1 Source 1: NYC-Airbnb-2023.csv

2.1.1 Location

- The dataset is obtained from Kaggle (<https://www.kaggle.com/datasets/godofoutcasts/new-york-city-airbnb-2023-public-data/data>)

2.1.2 Format: CSV (6.6 MB)

2.1.3 Important Variable Description

- 'neighborhood_group', 'neighborhood', 'latitude', 'longitude': containing the geographical information about each record.
- 'price': the amount needed to live in for one night.

- 'minimun_nights', ... 'availability_365': containing the information about the hardware of the room and the viewing statistics of the guests on the website.

2.1.4 Records Num & Time Period

- The dataset contains 42931 rows \times 18 columns, describing the Airbnb statistics in NYC for 2023.

2.2 Source 2: NYC_crime.csv

2.2.1 Location

- The dataset is obtained from Kaggle
(<https://www.kaggle.com/datasets/ajkarella/nyc-crime-stats>)

2.2.2 Format: CSV (688.2 MB)

- **Since the dataset is too large, in our zip file, we will only include the first 100 rows of crime.csv and NYC_crime.csv.**

2.2.3 Important Variable Description

- 'latitude', 'longitude': containing the geographical information about each record.
- 'pd_desc': Description of internal classification corresponding with PD code (more granular than Offense Description).
- 'ofns_desc': Description of offense corresponding with key code.
- 'law_cat_cd': Level of offense: felony, misdemeanor, violation.
- 'perp_race': Racial types of the criminals.

2.2.4 Records Num & Time Period

- The dataset contains 3881989 rows \times 18 columns, describing the crime statistics in NYC from 2006 to 2019.

2.3 Source 3: fullDownload.geojson

2.3.1 Location

- The dataset is obtained from
<https://dsl.richmond.edu/panorama/redlining/#loc=5/39.1/-94.58&text=downloads>

2.3.2 Format: GEOJSON (17.8MB)

2.3.3 Important Variable Description

- 'state', 'city', 'name': containing brief geographical information about the neighborhood in America.
- 'geometry', 'coordinates': containing detailed geometrical information about each neighborhood.
- 'holc_grade': Overall description of the dangerous extent of each neighborhood.

2.3.4 Records Num & Time Period

- The dataset contains 8878 records, describing the geographical statistics all over American neighborhoods.

3. Data Manipulation Methods

3.1 Source 1: NYC-airbnb-2023.csv

3.1.1 Initial Process

- Data was loaded from the NYC-Airbnb-2023.csv file. We first checked for data types and null values of each column. We dropped several columns that have little relevance to our main focus.

3.1.2 Locate each Airbnb record and add neighborhood_cd column

- See spec in 3.3.3.

Handling incorrect and missing values:

We first replace -1 in the neighborhood code column with NaN and drop these rows, since -1 means these records don't belong to any one of listed New York neighborhoods and will not contribute to our analysis. For the left regional missing values, they will automatically be omitted after we establish the new data frame.

3.2 Source 2: NYC_crime.csv

3.2.1 Initial Process

- Data was loaded from the nyc_crime.csv file. We first checked for data types and null values of each column. Then, we rename the last four ambiguous regional columns' titles to community district, borough boundaries, city council districts and police precincts.

3.2.2 Create new neighborhood-centered crime dataframe

- **Workflow:** New York city has 300+ neighborhoods, and we will assign each crime to its corresponding neighborhoods coded from 1-300+. Then, we will calculate crime statistics in each neighborhood and compress the 700MB giant dataset into 300+ rows. The new dataframe will be indexed by neighborhood code, followed by various columns statistically evaluating crime features in that neighborhood area. Finally, we will right merge the new dataframe with airbnb by the shared neighborhood code column.
- **Challenge1:** Locate each crime and add neighborhood_cd column
See spec in 3.3.3.

Handling incorrect and missing values:

We first replace -1 in the neighborhood code column with NaN and drop these rows, since -1 means these crimes don't belong to any one of listed New York neighborhoods and will not contribute to our analysis. For the left regional missing values, they will automatically be omitted after we establish the new data frame.

- **Challenge2:** Adding crime features for each coded neighborhood
column 1. annual crime rate

We first use datetime and add a year-wise column for each crime. Then, we group the data frame by neighborhood code and year, counting the number of crimes, unstacking the data frame and calculating the annual crime rate by mean method.

column 2-5. misdemeanor_rate, felony_rate, violation_rate, infraction_rate

Similar to 1, this time we group the data frame by neighborhood code and crime level including felony, misdemeanor, violation and infraction, calculating their annual happening frequency for each neighborhood.

column 6. dominant_ofns

We group the data frame by neighborhood code and specific crime category, counting the frequency, sorting the value to find out and keep the most frequent crime type in each neighborhood.

column 7. dominant_perp_race

Same as column 6, we seek to find out and keep the most frequent criminal's race in each neighborhood.

3.3 Source 3: fullDownload.geojson

3.3.1 Initial Process

- Data was loaded from the 'fullDownload.json'. Based on the features of the geojson file, we initialized the data frame using 'pd.read_json', and extracted the 'features' series. Then based on the 'state' and 'city'(see 2.3.3), we got the filtered series which only contains data in NYC.

3.3.2 Create class list

- We then defined a class called 'NYCDistrict' to store the geographical information of each district. Then we created a list called 'Districts' to store all the NYC districts. Now we had all the shape information of NYC districts.
- **Challenge:** However, we found that there are a great number of districts that don't have the name in the geojson file. We decided to use the order in the 'Districts' to distinguish the NYC district. The 'neighborhood_cd' column's process in the above two sources is based on this idea.

3.3.3 'neighborhood_cd' column process

- Based on the preexisting 'longitude' and 'latitude' value in 'NYC-airbnb-2023.csv', 'NYC_crime.csv', and the 'coordinates' information in 'NYCDistricts' which could form polygons, we assign the order of the district stored in 'Districts' to each record if the 'longitude' and 'latitude' of the record falls into the polygon. We assign -1 to records that don't fall into any polygon.

4. Analysis

After bringing together two different data resources, we are able to investigate further and find a new insight that could not have been answered with either data resource alone. Our analysis originated from three sections: EDA, a significant correlation between crime and house price and segmentation & clustering.

4.1 EDA(Exploratory Data Analysis)

4.1.1 Data Overview of Airbnb Crime

- In the process of data analysis, our data has merged both the crime and Airbnb datasets and conducted the initial cleaning of the missing value(See Part 3.2.2). In the Airbnb Crime dataset, values type consisted of int, object, and float. Values such as outliers(price, annual crime rate, etc) have been dropped by setting the parameters in the foreseeable figure or changing the percentile of the data to meet and grasp the critical aspects of the outcome.

4.1.2 Dominant Offenses & Price Range

- Given the Airbnb crime dataset, the visualization of all types of crime has been collected with respect to price segmentation. The 0-100 price range seems to have the highest overall count of offenses across all price ranges, suggesting a potential hotspot for various criminal activities. Higher price ranges (301-500 and above) have fewer total offenses, which could imply a correlation where more expensive areas have lower crime rates, or fewer listings in those price ranges are in high-crime areas. Specific offenses like “controlled substance, possession” and “marijuana, possession” are prevalent across multiple price ranges, which may reflect city-wide patterns in drug-related offenses.

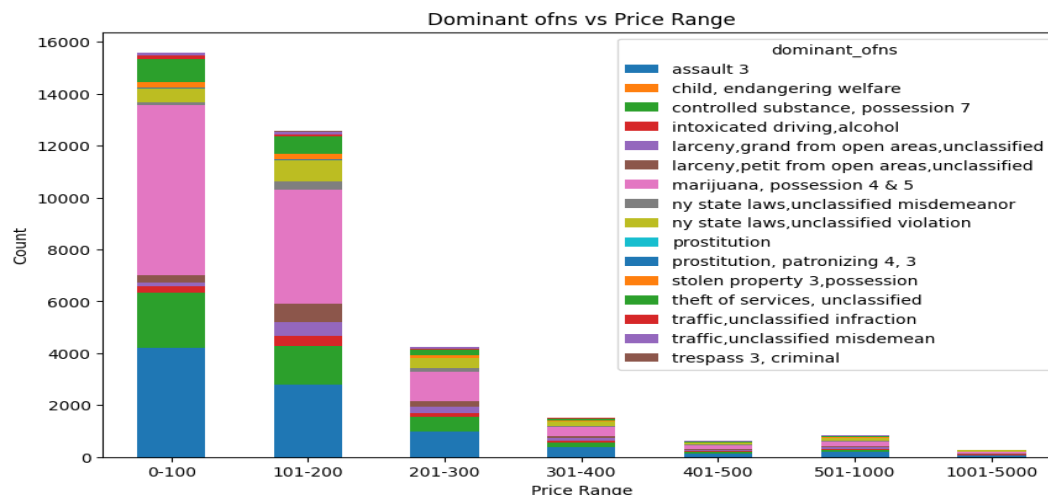


Fig4.1.2-1. Dominant ofns vs Price Range

4.2 Significant Correlation between Crime and House Price

4.2.1 Dominant Crime Type & Price

- Firstly, we want to explore whether the price distributions within each dominant crime type category will vary with each other. In order to better understand the significant difference, we first parse the airbnb_crime dataset into four subsequent datasets by room_type, the aim of which is to eliminate the prime influence by room types. However, since the number of items in shared room (465) and hotel room (107) is too small

compared with private room (15244) and entire room (19891), we will only show two box plots where x-axis is different dominant crimes and y-axis is the price:

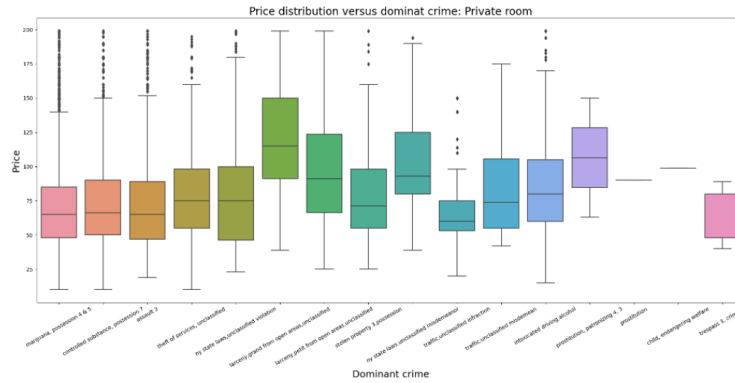


Fig4.2.1-1. Price distribution versus dominant crime: Private room

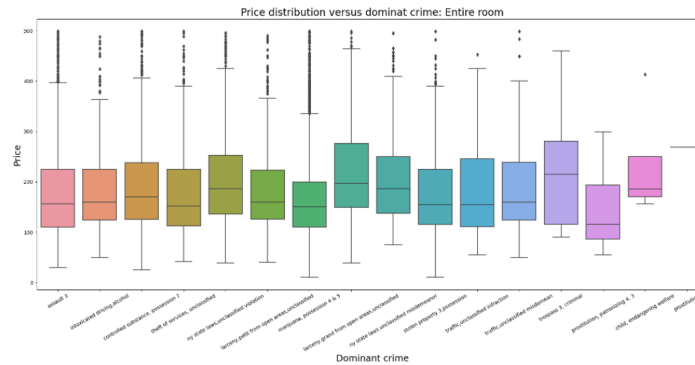


Fig4.2.1-2. Price distribution versus dominant crime: Entire room

According to Fig4.2.1-1 and Fig4.2.1-2, we can see that for some rooms lying in certain types of dominant crime neighborhoods, such as larceny and ny state law violation, the price distribution and their mean prices are evidently different from others. For further statistical proof, we conduct pairwise Tukeyhsd analysis as follows (Since the output is too long, we will present a part of the screenshot. The whole output can be referred in the combined_analysis.ipynb):

```
res2 = pairwise_tukeyhsd(private_room['price'], private_room['dominant_crime'])
res2.summary()
```

| group1 | group2 | meandiff | p-value | lower | upper | reject |
|----------------------------|--|----------|---------|-----------|----------|--------|
| assault 3 | child, endangering welfare | 19.7582 | 1.0 | -108.1272 | 167.5995 | False |
| assault 3 | controlled substance, possession 7 | -1.4488 | 0.9997 | -6.2926 | 3.3405 | False |
| assault 3 | intoxicated driving, alcohol | 30.3327 | 0.0 | 13.7296 | 36.9347 | True |
| assault 3 | larceny, grand from open areas, unclassified | 76.3815 | 0.0 | 66.379 | 86.384 | True |
| assault 3 | larceny, petit from open areas, unclassified | 107.078 | 0.0 | 26.1035 | 48.0302 | True |
| assault 3 | marijuana, possession & a.s | -3.4301 | 0.9997 | -10.777 | 3.2169 | False |
| assault 3 | ny state law, unclassified misdemeanor | 43.407 | 0.0 | 28.9468 | 57.9091 | True |
| assault 3 | ny state law, unclassified violation | 20.1674 | 0.0 | 12.5119 | 28.6229 | True |
| assault 3 | prostitution | 50.7992 | 1.0 | -115.5483 | 165.2206 | False |
| assault 3 | prostitution, patronizing & s | 27.2362 | 1.0 | -88.5483 | 163.2206 | False |
| assault 3 | stolen property, possession | 12.9436 | 0.9995 | -6.3357 | 26.2229 | False |
| assault 3 | theft of services, unclassified | 13.8269 | 0.0 | 7.2923 | 20.3595 | True |
| assault 3 | traffic, unclassified infraction | -8.2986 | 0.8795 | -26.1251 | 7.5278 | False |
| assault 3 | traffic, unclassified misdemeanor | 33.7688 | 0.0 | 15.4766 | 62.0369 | True |
| assault 3 | weapons, a.s, criminal | -19.2638 | 1.0 | -87.8475 | 61.3198 | False |
| child, endangering welfare | controlled substance, possession 7 | -21.1922 | 1.0 | -109.6746 | 166.6902 | False |
| child, endangering welfare | intoxicated driving, alcohol | 5.896 | 1.0 | -172.8986 | 184.7918 | False |
| child, endangering welfare | larceny, grand from open areas, unclassified | 84.6453 | 0.0005 | -12.4435 | 124.7442 | False |
| child, endangering welfare | larceny, petit from open areas, unclassified | 17.3357 | 1.0 | -180.8188 | 195.4901 | False |
| child, endangering welfare | marijuana, possession & a.s | -23.9683 | 1.0 | -207.0214 | 154.6848 | False |
| child, endangering welfare | ny state law, unclassified misdemeanor | 23.6708 | 1.0 | -164.7232 | 212.0658 | False |
| child, endangering welfare | ny state law, unclassified violation | 0.8312 | 1.0 | -177.688 | 179.3513 | False |
| child, endangering welfare | prostitution | -8.0 | 1.0 | -228.8093 | 208.8093 | False |
| child, endangering welfare | prostitution, patronizing & s | 1.5 | 1.0 | -210.3093 | 223.3093 | False |
| child, endangering welfare | stolen property, possession | -6.7926 | 1.0 | -185.1054 | 171.6203 | False |
| child, endangering welfare | theft of services, unclassified | -5.8997 | 1.0 | -183.8249 | 172.0455 | False |

Fig4.2.2-1 Multiple Comparison of Means - Tukey HSD (Private room)

```
res2 = pairwise_tukeyhsd(entire_room['price'], entire_room['dominant_crime'])
res2.summary()
```

| group1 | group2 | meandiff | p-value | lower | upper | reject |
|----------------------------|--|----------|---------|-----------|----------|--------|
| assault 3 | child, endangering welfare | 10.4515 | 0.9961 | -62.6516 | 203.5552 | False |
| assault 3 | controlled substance, possession 7 | 10.8086 | 0.0003 | 2.8542 | 18.7629 | True |
| assault 3 | intoxicated driving, alcohol | 4.9105 | 0.9997 | -8.2975 | 18.0395 | False |
| assault 3 | larceny, grand from open areas, unclassified | 80.0769 | 0.0 | 36.079 | 104.0771 | True |
| assault 3 | larceny, petit from open areas, unclassified | 5.1482 | 0.9439 | -5.1475 | 15.4388 | False |
| assault 3 | marijuana, possession & a.s | -12.4884 | 0.0 | -18.2847 | -6.6721 | True |
| assault 3 | ny state law, unclassified misdemeanor | 23.0203 | 0.0 | 32.2308 | 41.8238 | True |
| assault 3 | ny state law, unclassified violation | 26.7425 | 0.0 | 17.7206 | 35.7644 | True |
| assault 3 | prostitution | 85.8313 | 0.0087 | -208.6548 | 380.5374 | False |
| assault 3 | prostitution, patronizing & s | -17.7802 | 0.9362 | -119.9873 | 84.4448 | False |
| assault 3 | stolen property, possession | 0.1786 | 1.0 | -15.347 | 15.7051 | False |
| assault 3 | theft of services, unclassified | -0.529 | 1.0 | -11.2275 | 10.1694 | False |
| assault 3 | traffic, unclassified infraction | -0.9764 | 1.0 | -24.8719 | 22.9191 | False |
| assault 3 | traffic, unclassified misdemeanor | 10.7213 | 0.9429 | -10.6791 | 32.1217 | False |
| assault 3 | weapons, a.s, criminal | -0.9391 | 0.9875 | -40.7602 | 37.8826 | False |
| child, endangering welfare | controlled substance, possession 7 | -44.9427 | 0.0097 | -103.1105 | 13.225 | False |
| child, endangering welfare | intoxicated driving, alcohol | -35.3378 | 0.0088 | -109.0306 | 38.0101 | False |
| child, endangering welfare | larceny, grand from open areas, unclassified | -17.3301 | 1.0 | -168.8323 | 131.0813 | False |
| child, endangering welfare | larceny, petit from open areas, unclassified | -50.3051 | 0.9998 | -108.6327 | 88.0224 | False |
| child, endangering welfare | marijuana, possession & a.s | -47.8197 | 0.0714 | -219.0201 | 81.1818 | False |
| child, endangering welfare | ny state law, unclassified misdemeanor | -28.450 | 1.0 | -177.1338 | 120.2838 | False |
| child, endangering welfare | ny state law, unclassified violation | -18.7077 | 1.0 | -176.9525 | 118.537 | False |
| child, endangering welfare | prostitution | 14.8 | 1.0 | -280.0796 | 360.0796 | False |
| child, endangering welfare | prostitution, patronizing & s | -83.1875 | 0.8827 | -284.8073 | 72.3223 | False |
| child, endangering welfare | stolen property, possession | -45.2727 | 0.0088 | -204.0549 | 93.5095 | False |
| child, endangering welfare | theft of services, unclassified | -45.9903 | 0.9902 | -204.2428 | 92.2621 | False |
| child, endangering welfare | traffic, unclassified infraction | -56.4277 | 0.0963 | -206.3146 | 93.4593 | False |
| child, endangering welfare | traffic, unclassified misdemeanor | -44.73 | 0.9997 | -194.2595 | 104.7975 | False |

Fig4.2.2-2 Multiple Comparison of Means - Tukey HSD (Entire room)

- By Tukey HSD, from entire room and private room groups, room prices lying in larceny and ny state law violation dominant crime neighborhoods tend to be obviously significantly different from room prices lying in other types of dominant crime neighborhoods. It is because their pairwise p-values are lower than the chosen significance level (commonly 0.05), which suggests the difference in means between the two groups is statistically significant.
- For our understanding, a possible explanation can be given that these neighborhoods' airbnb room prices are normally set higher than other neighborhoods, which means the owner or customers who choose these rooms could be richer. As a result, the house owners or customers are easier to be chosen as a larceny target.

4.2.2 Annual Crime Rate & Price

- Secondly, aside from categorical crime data, we also look into quantitative features of annual crime rate as continuous variables in detail. We want to answer whether these additional features can help us explain the price better. Thus, we perform the OLS (ordinary least squared) model where room_type and four level annual_crime_rate (felony, misdemeanor, violation and infraction) as independent variables and room price as dependent variable:

| OLS Regression Results | | | | | | |
|--------------------------------|------------------|---------------------|-------------|-------|----------|---------|
| ===== | | | | | | |
| Dep. Variable: | Q('price') | R-squared: | 0.312 | | | |
| Model: | OLS | Adj. R-squared: | 0.312 | | | |
| Method: | Least Squares | F-statistic: | 2233. | | | |
| Date: | Wed, 06 Dec 2023 | Prob (F-statistic): | 0.00 | | | |
| Time: | 02:59:46 | Log-Likelihood: | -1.9715e+05 | | | |
| No. Observations: | 34448 | AIC: | 3.943e+05 | | | |
| Df Residuals: | 34440 | BIC: | 3.944e+05 | | | |
| Df Model: | 7 | | | | | |
| Covariance Type: | nonrobust | | | | | |
| ===== | | | | | | |
| | coef | std err | t | P> t | [0.025 | 0.975] |
| Intercept | 185.4820 | 0.689 | 269.083 | 0.000 | 184.131 | 186.833 |
| Q('room_type')[T.Hotel room] | 26.4795 | 7.579 | 3.494 | 0.000 | 11.624 | 41.335 |
| Q('room_type')[T.Private room] | -97.3435 | 0.812 | -119.831 | 0.000 | -98.936 | -95.751 |
| Q('room_type')[T.Shared room] | -103.5223 | 3.519 | -29.416 | 0.000 | -110.420 | -96.624 |
| Q('annual_felony_rate') | -0.0074 | 0.003 | -2.182 | 0.029 | -0.014 | -0.001 |
| Q('annual_misdemeanor_rate') | -0.0085 | 0.001 | -5.799 | 0.000 | -0.011 | -0.006 |
| Q('annual_violation_rate') | 0.1225 | 0.006 | 19.944 | 0.000 | 0.110 | 0.135 |
| Q('annual_infraction_rate') | -0.9852 | 0.088 | -11.160 | 0.000 | -1.158 | -0.812 |
| ===== | | | | | | |
| Omnibus: | 9704.058 | Durbin-Watson: | 1.885 | | | |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 26872.036 | | | |
| Skew: | 1.503 | Prob(JB): | 0.00 | | | |
| Kurtosis: | 6.113 | Cond. No. | 4.09e+04 | | | |
| ===== | | | | | | |

Fig4.2.2-1 OLS Regression Results

- According to the OLS regression result:
1. The model indicates that annual crime rates are significant predictors (p-value < 0.05) of Airbnb prices. The negative coefficients for annual felony, misdemeanor, and infraction rates suggest that higher crime rates in these categories are associated with lower prices.

- Despite the statistical significance of the predictors, the moderate R-squared value implies that other factors not included in the model also play a significant role in determining Airbnb prices, which means the regression power of annual crime rates is very limited. Thus, we have some reasons to believe that room prices are less likely to be determined by features of crime rates.

4.3 Segmentation & Clustering

4.3.1 Clustering (exclude price)

- Firstly, we want to find the potential patterns inside the records without the influence of price. Thus we dropped the 'id' and 'price' columns. Then we used the pipeline to first scale each remaining column, then reduce the dimension to 5, and finally use K-means to cluster the records.
- We plotted the cluster visualization and the average silhouette score with the K-means clustering number varying from 2 to 20. We found when the cluster number is 7, the average silhouette score is the highest (0.318).

Silhouette analysis for KMeans clustering on sample data with n_clusters = 7

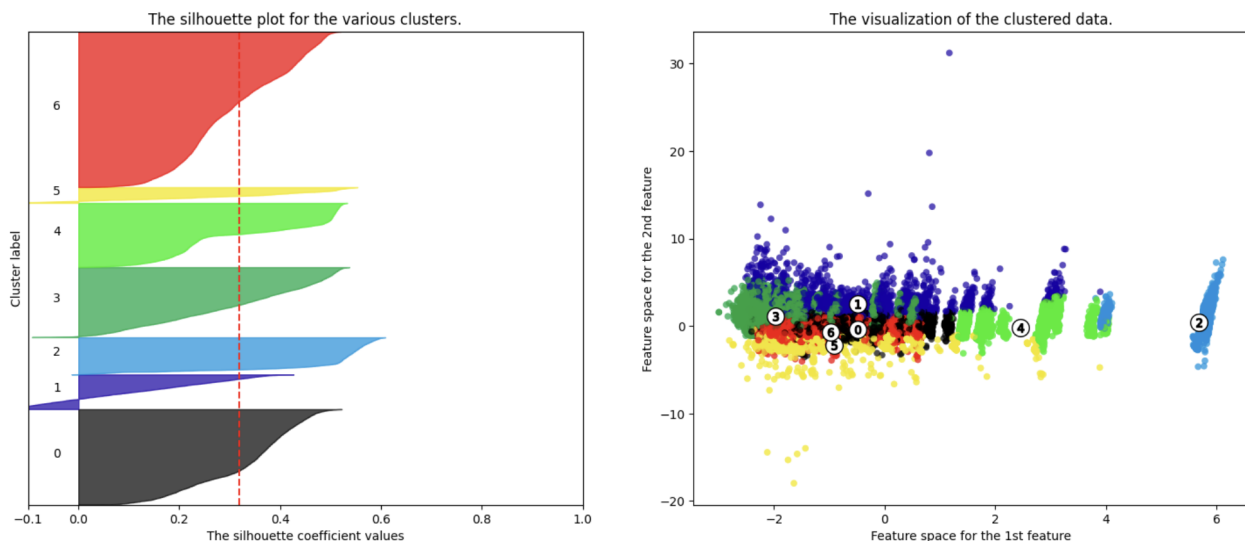


Fig4.3.1 Optimal Clustering Result

4.3.2 Clustering Result Analysis

- Based on the optimal clustering result, we assign the segments to the original data (including 'id' and 'price'). An intuitive idea is that each cluster's price distribution should be in line with each crime rate if there exist certain relationships. We decided to compare them using box plots.

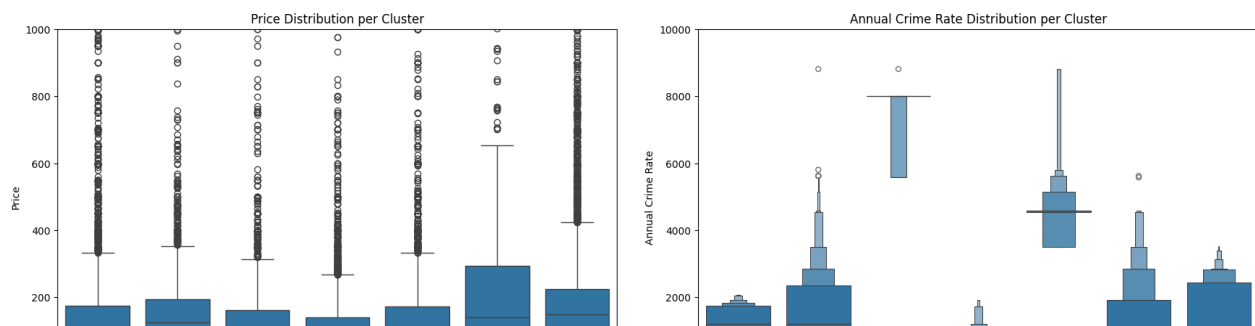


Fig4.3.2 Price/Annual Crime Rate per Cluster

- From the above two boxplots, we cannot find a direct relationship between the clusters and the crime rate. For example, Cluster 2 has the highest crime rate while Cluster 3 has the lowest crime rate, but there is no significant difference in prices between the two clusters.

4.3.3 Comparison of Visualization (clusters exclude price, clusters only with price, clusters from geojson)

- Now we attempt to visualize the clusters on the map. We first assign the neighborhood to the most frequent cluster in the neighborhood. Then we plot the clusters on the map. We decided to have three map visualizations based on clusters excluding price, clusters only with price, and clusters from geojson.

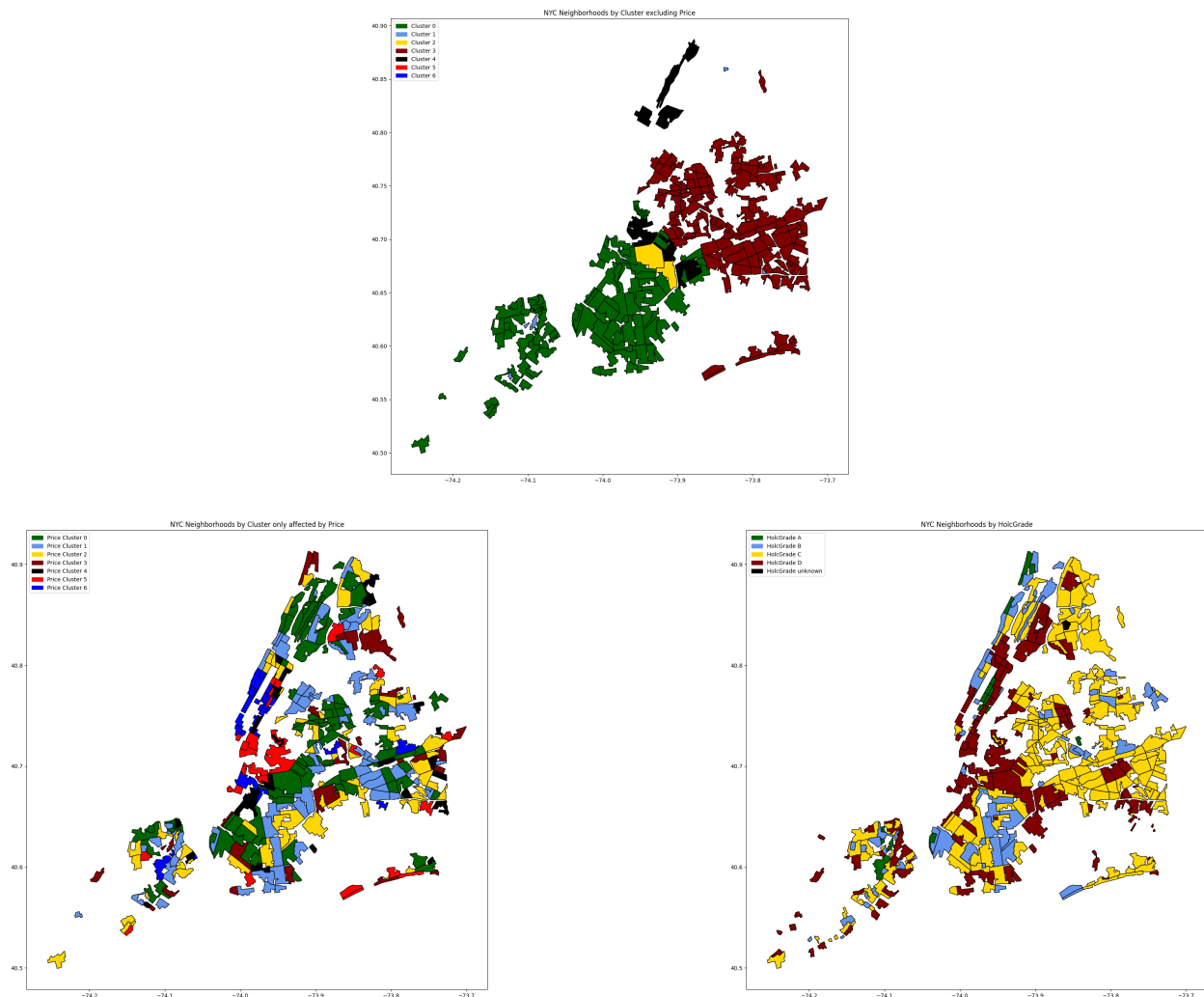


Fig4.3.3 Map Visualizations (top: clusters without price; left bottom: clusters by price; right bottom: clusters by 3.3)

- We find that the clusters from K-means do not correlate with the price, but are determined only by the region of the NYC city. Additionally, the general dangerous extent has no obvious correlations with price from the comparison of the bottom two. Thus, the crime rate may not be toppest main factor affecting the price of Airbnb.

5. Conclusion

Our study delved into how crime rates and Airbnb pricing interact within NYC's neighborhoods, revealing that prices aren't solely driven by crime but also by a tapestry of regional traits. Our spatial analysis with K-means clustering, Tukey HSD and geographic visualization, offer a guide for the Airbnb community and a strategic tool for urban development, integrating economic and safety considerations.

5.1 Limitation

While our analysis has yielded significant observations, the dataset still has some flaws. One notable constraint is the reliance on available data, which may not fully capture the real-time fluctuations in crime and economic factors. Furthermore, our clustering analysis indicated that the correlation between Airbnb pricing and crime is not as robust as initially hypothesized, underscoring the need for a more granular approach that considers additional variables and market forces.

Our visualizations, though comprehensive, also underline the absence of a clear-cut relationship between the perceived danger in various regions and the pricing of Airbnb listings. This highlights the complexity of the market and suggests that additional factors must be accounted for to accurately predict and strategize pricing models in correlation with crime statistics.

5.2 Next Step

As we look to the future, we aim to expand our dataset, incorporate real-time analytics, and explore a wider array of contributing factors. This will not only enhance the precision of our insights but also provide a more actionable framework for all stakeholders involved.

6. Statement of Work

6.1 Contribution Report

6.1.1 Our team collaborated on the project based on the following aspects:

| | |
|--------------|--|
| Keye Chen | Data sources, geojson file data manipulation, clustering analysis |
| Fengyu Zhang | Data manipulation of aggregation and merge, correlation analysis |
| Boyan Wu | Motivation, EDA, Data Visualization, Conclusion and Statement of Work, |

Assessment and Improvement:

Our team has distributed work and divided it to leverage individual strengths. Each member contributes to a different facet of the project. However, we can improve our collaboration in listed below in the future:

Regular Check-Ins: Schedule periodic meetings to ensure everyone is on track and to discuss any challenges or findings, making it easier for team members to collaborate and pick up where others left off.

Cross Reviews: Occasionally, swap roles or review each other's work would reduce the risk of siloed information, and increase team members' skills across different areas.

Feedback Improvement: Each team member can raise the feedback on the overall performance to reflect on continuous improvement in team dynamics and project outcomes.