# MSDS 593 Final Project

## AIRBNB OPEN DATA

Bhumika Srinivas

Eren Bardak

Sissi Shen

# INSIGHTFUL DATA EXPLORATION

• Dropping duplicate rows and unnecessary columns

```python
#DROPPING DUPLICATE ROWS
df = df.drop_duplicates()

#DROPPING COLUMNS : country, country code, license, host name, host id
df = df.drop(['country','country code','license','host name'], axis=1)
```

• Corrected misspellings and Imputed missing values

```python
#Merging the misspellings: manhatan and brookln
#COLUMN : Neighbourhood Group
df['neighbourhood group'] = df['neighbourhood group'].replace('manhatan','Manhattan')
df['neighbourhood group'] = df['neighbourhood group'].replace('brookln','Brooklyn')

#Replaces the missing value with the average price
df['service fee'] = df['service fee'].fillna(df['service fee'].mean())

#Imputting missing values with the mean lat long of their neighborhoods

#Manhattan East Village
manhattan_eastvillage = df[(df['neighbourhood group'] == 'Manhattan') & (df['neighbourhood'] == 'East Village')]
df.loc[(df['lat'].isna()) & (df['id'] == '1442624'), 'lat'] = manhattan_eastvillage['lat'].mean()
df.loc[(df['long'].isna()) & (df['id'] == '1442624'), 'long'] = manhattan_eastvillage['long'].mean()
#Manhattan West Village
manhattan_westvillage = df[(df['neighbourhood group'] == 'Manhattan') & (df['neighbourhood'] == 'West Village')]
df.loc[(df['lat'].isna()) & (df['id'] == '1450908'), 'lat'] = manhattan_westvillage['lat'].mean()
df.loc[(df['long'].isna()) & (df['id'] == '1450908'), 'long'] = manhattan_westvillage['long'].mean()
```

• Change of data types

```python
#Changing the type as category

#COLUMN : id
df['id'] = df['id'].astype('category')
#COLUMN : Host Identity Verified
df['host_identity_verified'] = df['host_identity_verified'].astype('category')
#COLUMN : Neighbourhood Group
df['neighbourhood group'] = df['neighbourhood group'].astype('category')
#COLUMN : Neighbourhood
df['neighbourhood'] = df['neighbourhood'].astype('category')
#COLUMN : Instant Bookable
df['instant_bookable'] = df['instant_bookable'].astype('category')
```
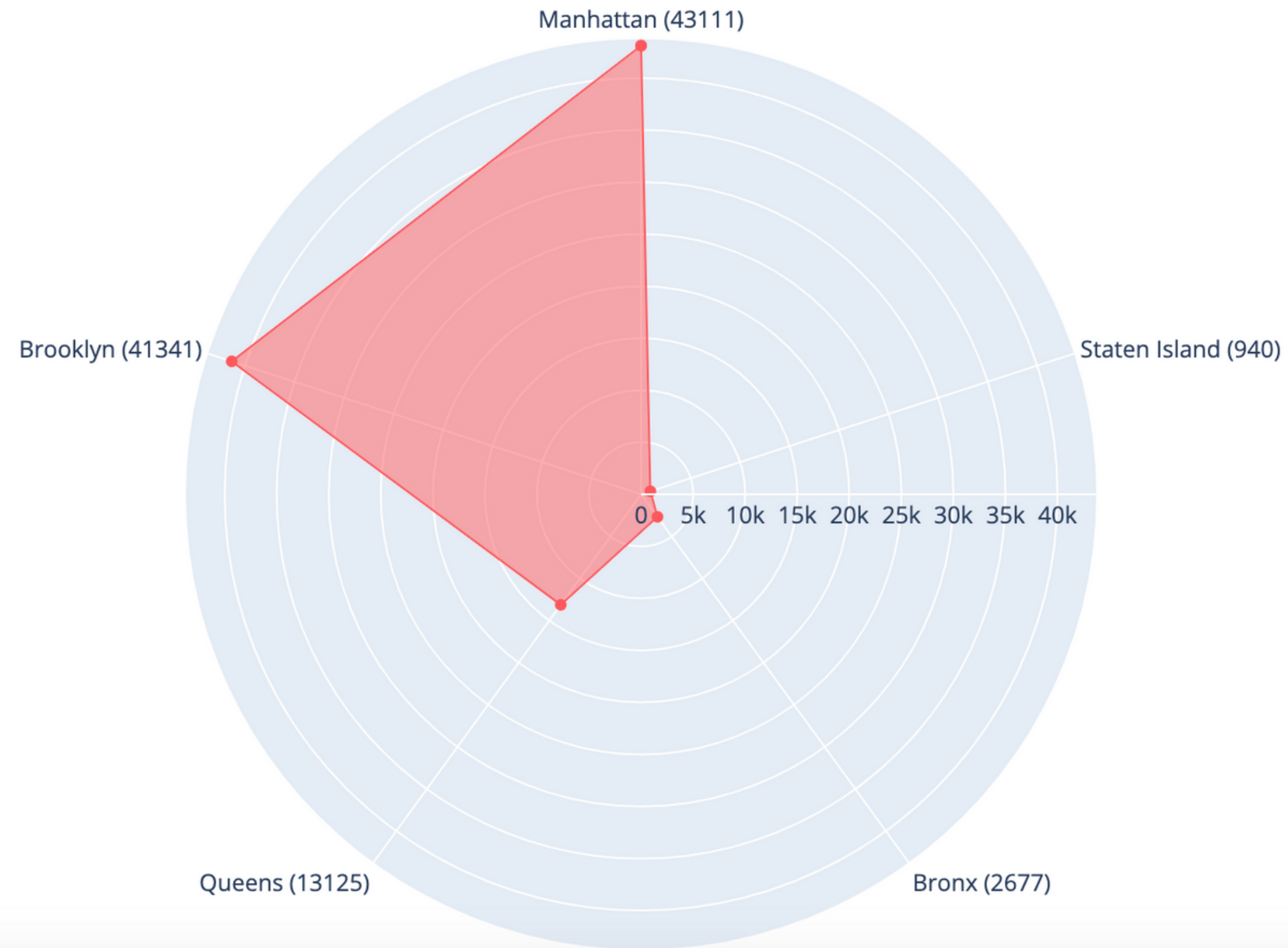
# STATISTICAL SNAPSHOTS

**26 Variables │ 102599 Records** ⟶ **15 Variables │ 101221 Records**

| | number of reviews | price | review rate number | calculated host listings count |
|---|---|---|---|---|
| **count** | 101221.000000 | 101221.000000 | 101221.000000 | 101221.000000 |
| **mean** | 27.476502 | 625.469189 | 3.279468 | 8.059617 |
| **std** | 49.441459 | 331.186977 | 1.283043 | 32.296895 |
| **min** | 0.000000 | 50.000000 | 1.000000 | 1.000000 |
| **25%** | 1.000000 | 341.000000 | 2.000000 | 1.000000 |
| **50%** | 7.000000 | 625.355580 | 3.000000 | 1.000000 |
| **75%** | 30.000000 | 912.000000 | 4.000000 | 2.000000 |
| **max** | 1024.000000 | 1200.000000 | 5.000000 | 332.000000 |

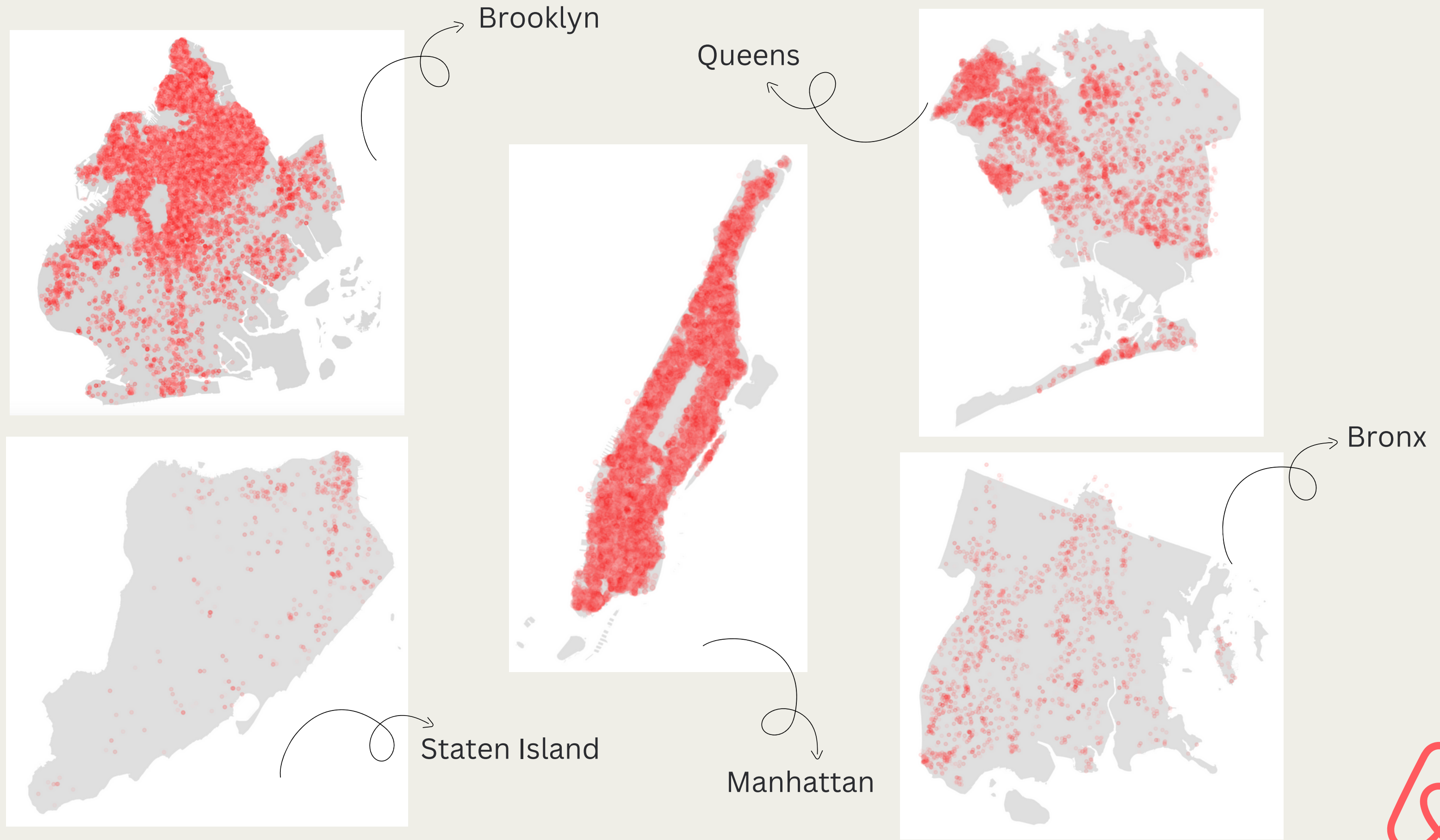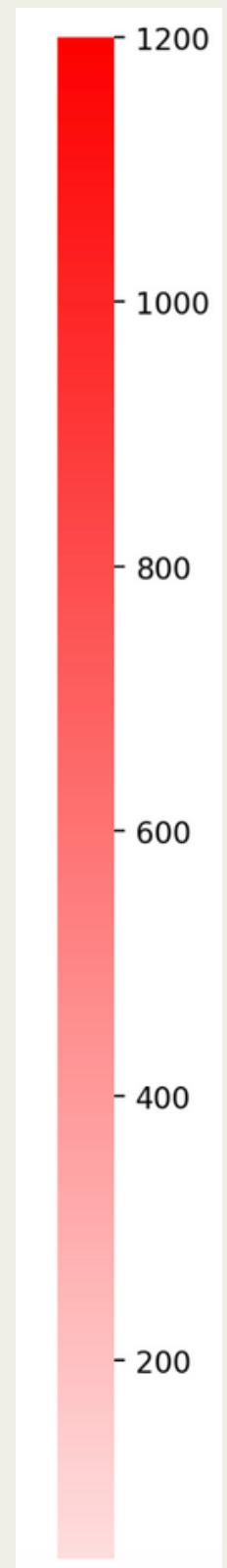# #1 Neighbourhood Showdown: *Listing Counts Across Groups*



Unveiling Listing Counts Across Neighbourhood Groups

Manhattan (43111)
Staten Island (940)
Brooklyn (41341)
Bronx (2677)
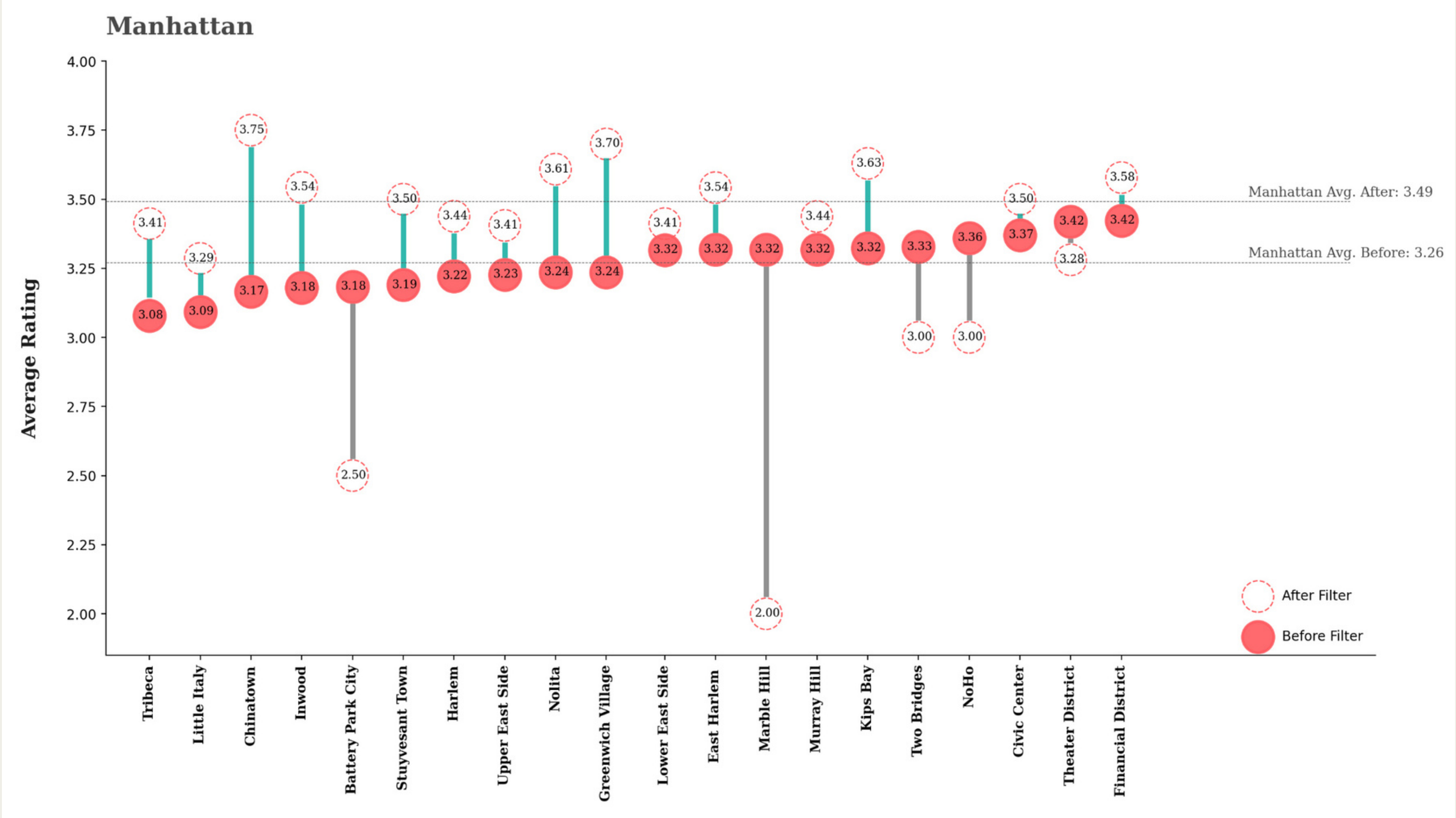Queens (13125)

0  5k  10k  15k  20k  25k  30k  35k  40k

# #2 Luxury Landscape: *Price Ranges by Neighbourhood Groups*

Price ($)

1200

1000

800

600

400

200

Brooklyn

Queens

Bronx

Staten Island

Manhattan

# #3 Enhanced Manhattan Ratings: *Excluding Old, Long-Term, and Inactive Listings*



Manhattan

Average Rating

Manhattan Avg. After: 3.49
Manhattan Avg. Before: 3.26

| Neighborhood | After Filter | Before Filter |
|---|---|---|
| Tribeca | 3.41 | 3.08 |
| Little Italy | 3.29 | 3.09 |
| Chinatown | 3.75 | 3.17 |
| Inwood | 3.54 | 3.18 |
| Battery Park City | 2.50 | 3.18 |
| Stuyvesant Town | 3.50 | 3.19 |
| Harlem | 3.44 | 3.22 |
| Upper East Side | 3.41 | 3.23 |
| Nolita | 3.61 | 3.24 |
| Greenwich Village | 3.70 | 3.24 |
| Lower East Side | 3.41 | 3.32 |
| East Harlem | 3.54 | 3.32 |
| Marble Hill | 2.00 | 3.32 |
| Murray Hill | 3.44 | 3.32 |
| Kips Bay | 3.63 | 3.32 |
| Two Bridges | 3.00 | 3.33 |
| NoHo | 3.00 | 3.36 |
| Civic Center | 3.50 | 3.37 |
| Theater District | 3.28 | 3.42 |
| Financial District | 3.58 | 3.42 |

- - - - After Filter
● Before Filter

# #3 Enhanced Manhattan Ratings: *Excluding Old, Long-Term, and Inactive Listings*



**Manhattan**

Average Rating

Manhattan Avg. After: 3.49
Manhattan Avg. Before: 3.26

| District | Before | After |
|---|---|---|
| Tribeca | 3.08 | 3.41 |
| Little Italy | 3.09 | 3.29 |
| Chinatown | 3.17 | 3.75 |
| Inwood | 3.18 | 3.54 |
| Battery Park City | 3.18 | 2.50 |
| Stuyvesant Town | 3.19 | 3.50 |
| Harlem | 3.22 | 3.44 |
| Upper East Side | 3.23 | 3.41 |
| Nolita | 3.24 | 3.61 |
| Greenwich Village | 3.24 | 3.70 |
| Lower East Side | 3.32 | 3.41 |
| East Harlem | 3.32 | 3.54 |
| Marble Hill | 3.32 | 2.00 |
| Murray Hill | 3.32 | 3.44 |
| Kips Bay | 3.32 | 3.63 |
| Two Bridges | 3.33 | 3.00 |
| NoHo | 3.36 | 3.00 |
| Civic Center | 3.37 | 3.50 |
| Theater District | 3.42 | 3.28 |
| Financial District | 3.42 | 3.58 |

Less than 10 listings

After Filter
Before Filter

# #4 Neighbourhood Gems: *Top Rated Hotspots Within Each Group*



Neighborhood popularity

Average number of reviews

Top 5 frequent neighbourhoods

**Bronx** — Below Average, Above Average
Fordham, Kingsbridge, Longwood, Mott Haven, Wakefield
Bronx Average

**Manhattan**
East Village, Harlem, Hell's Kitchen, Upper East Side, Upper West Side
Manhattan Average

**Brooklyn**
Bedford-Stuyvesant, Bushwick, Crown Heights, Greenpoint, Williamsburg
Brooklyn Average

**Queens**
Astoria, Flushing, Long Island City, Ridgewood, Sunnyside
Queens Average

**Staten Island**
Concord, St. George, Stapleton, Tompkinsville, West Brighton
Staten Island Average

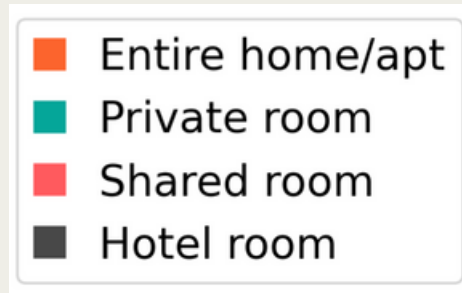# #5 Room Type: *Proportions Across Neighbourhood Groups*



Brooklyn

Queens

Manhattan

Staten Island

Bronx

Legend:
- Entire home/apt
- Private room
- Shared room
- Hotel room

# #5 Room Type: *Proportions Across Neighbourhood Groups*



*Brooklyn*

*Queens*

Older Than 10 Years

**Legend:**
- ■ Entire home/apt
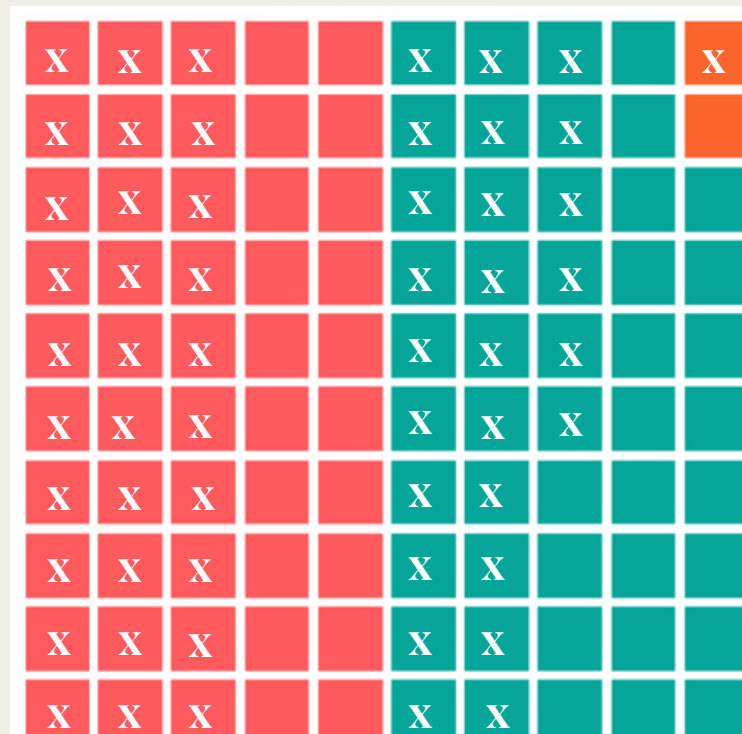- ■ Private room
- ■ Shared room
- ■ Hotel room

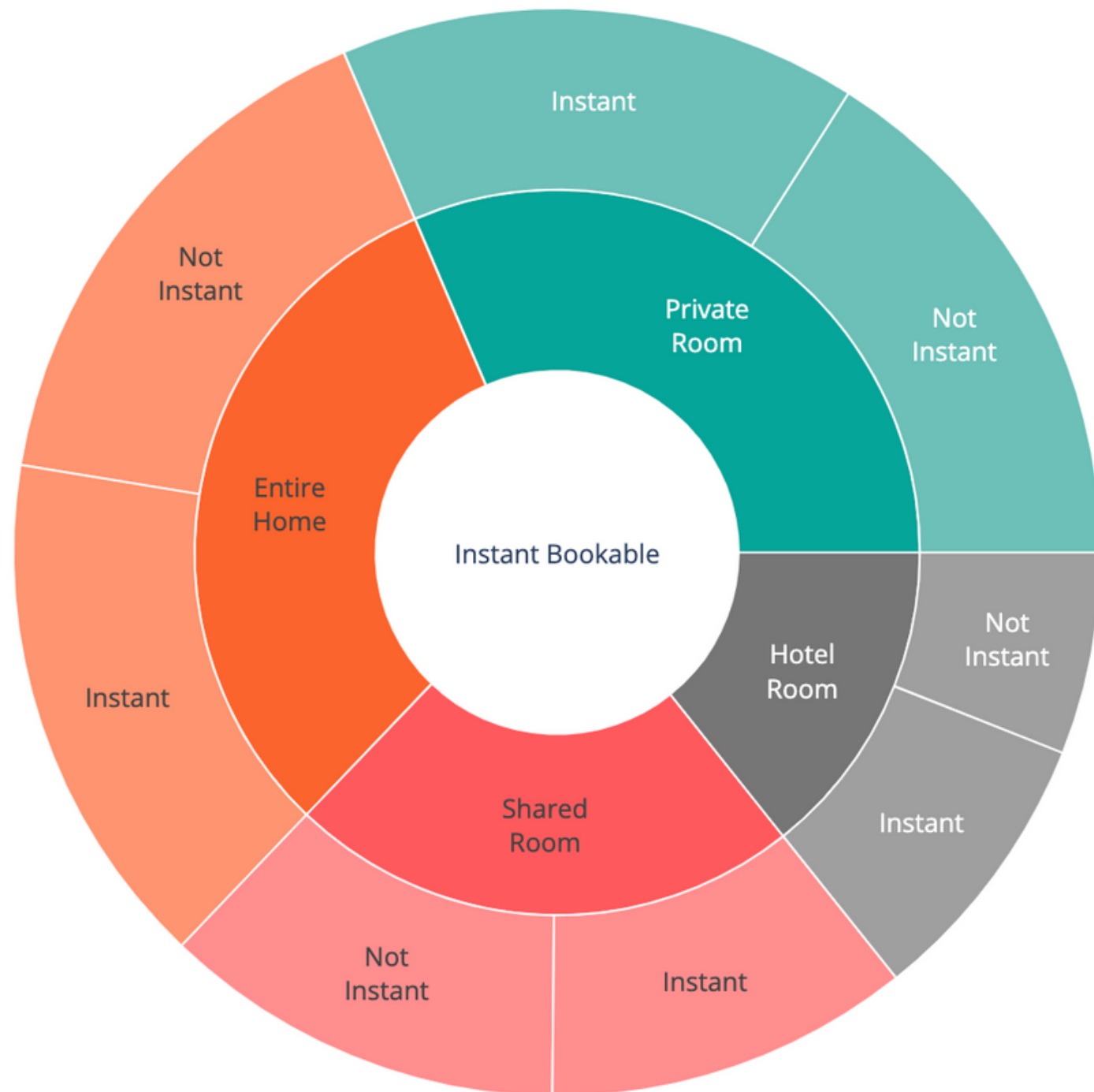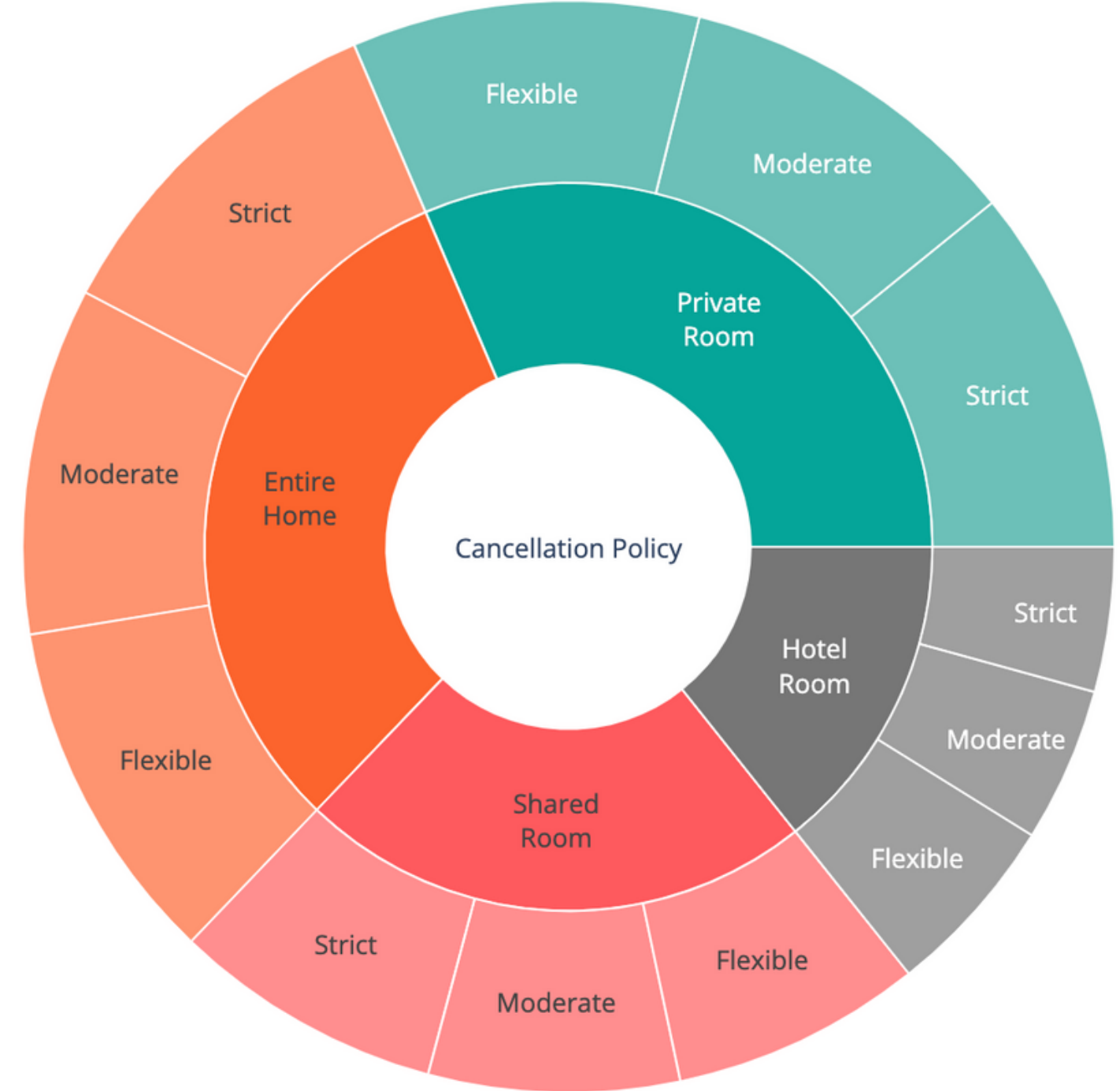*Manhattan*

*Staten Island*

*Bronx*

# #5 Room Type: *Proportions Across Neighbourhood Groups*

## 1 Instant Bookable



## 2 Cancellation Policy

# #6 Host Analysis: *Review Scores and Customer Engagement*



**Average Review Score and Number of Reviews by Host Types**

Hotel Hosts
Totel Number of Hosts: ~10
Review Count: 3.47
Review Score: 3.4678

Professional Small
Property Managers
Totel Number of Hosts: ~10871
Review Count: 36.60
Review Score: 3.3253

Individual Hosts
Totel Number of Hosts: ~62591
Review Count: 24.68
Review Score: 3.2463

Professional Large Property Managers
Totel Number of Hosts: ~84
Review Count: 7.71
Review Score: 3.3294

Average review score per lisitng

Average number of reviews per listing