

# DEVELOPMENT TEAM PROJECT: AIRBNB NEW YORK CITY 2019

## ANALYTIC REPORT

Maria Ingold, Murthy Kanuri, Khoi Dang (Danty), Ahmed Husain  
Unit 6  
Machine Learning  
University of Essex Online  
2 December 2024

## **CONTENTS**

INTRODUCTION.....	3
RATIONALE.....	4
DATA UNDERSTANDING AND PREPARATION.....	6
DATA LEARNINGS.....	8
Price .....	8
Location .....	9
Room Type .....	11
Host Characteristics .....	12
Correlation .....	14
Clustering .....	15
INSIGHTS .....	18
CONCLUSION .....	19
REFERENCES.....	20
APPENDIX.....	23

## INTRODUCTION

Founded in 2008, Airbnb has generated \$250 billion revenue from five million hosts (Airbnb, N.D.a; Ding et al., 2023). Despite revenue increases (Figure 1), year-over-year (Y/Y) growth declined below pre-COVID levels, prompting further Artificial Intelligence (AI) research (Airbnb, N.D.b).

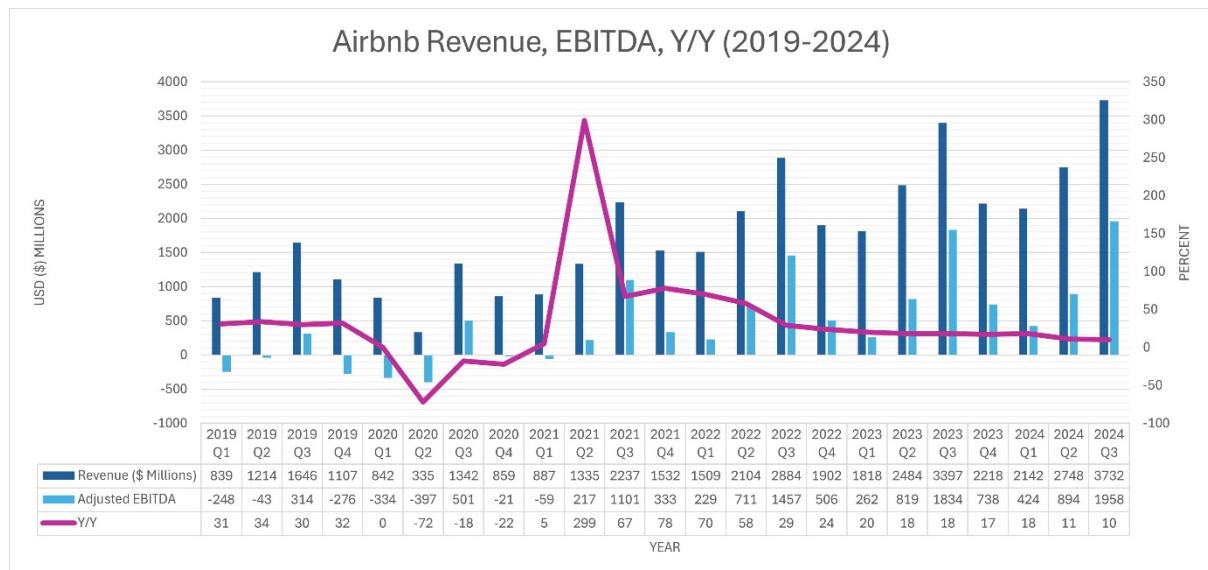
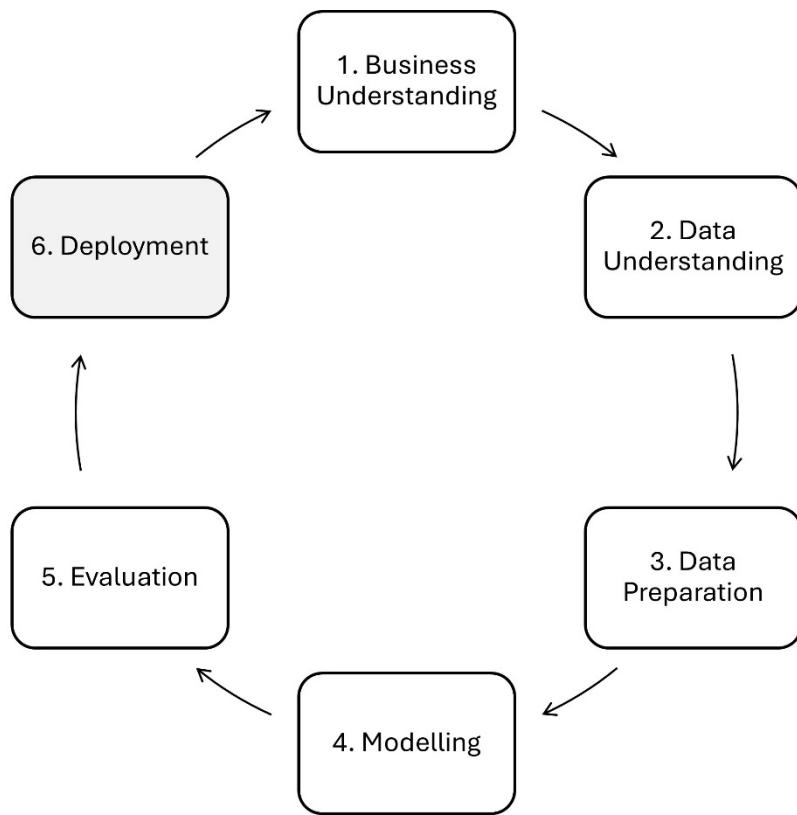


FIGURE 1 | Airbnb

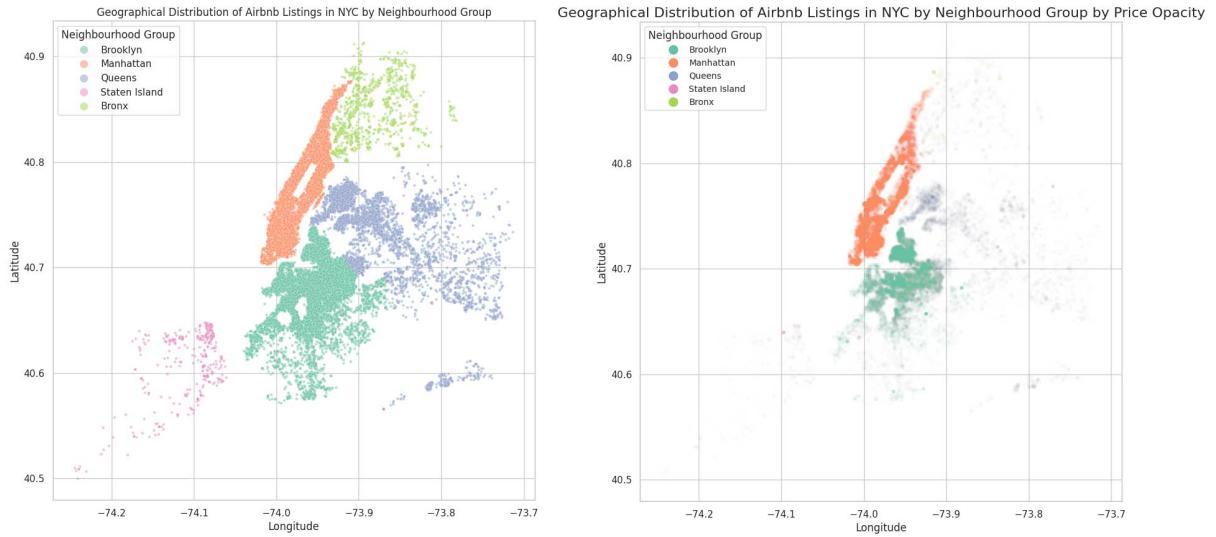
This report applies Cross Industry Standard Process for Data Mining (CRISP-DM) (Figure 2) to the 2019 New York City (NYC) Airbnb dataset (Niakšu, 2015; Kaggle, 2021). Through exploratory data analysis (EDA) and unsupervised machine learning (ML), it visualises trends and generates pricing recommendations to optimise revenue.



**FIGURE 2 | CRISP-DM**

## RATIONALE

Quarterly year-over-year revenue growth (Figure 1) reflects seasonal trends and company health (Taylor & Almeida, 2022). NYC's price-inelasticity enables raising prices to increase revenue, however, market dynamics need consideration (Gunter et al., 2020; Akalın & Alptekin, 2024). Therefore, this study addresses: "*How can Airbnb NYC optimise revenue growth through tailored pricing recommendations based on location, room type, and host characteristics?*"



**FIGURE 3 |** Airbnb listing neighbourhoods and highest price densities

Location affects pricing (Figure 3), including proximity to amenities (Akalin & Alptekin, 2024). Room type influences market positioning, impacting competition and regulation (Gunter et al., 2020). Host characteristics like listing number, availability, and reviews potentially reveal professional hosts targeted for regulation (Table 1).

**TABLE 1 |** Potential price influencers

<b>Location</b>	<ul style="list-style-type: none"> <li>• General: neighbourhood_group</li> <li>• Specific: neighbourhood, latitude, longitude</li> </ul>
<b>Room Type</b>	<ul style="list-style-type: none"> <li>• room_type: private, shared, entire</li> </ul>
<b>Host Characteristics</b>	<ul style="list-style-type: none"> <li>• calculated_host_listings_count</li> <li>• minimum_nights</li> <li>• availability_365</li> <li>• number_of_reviews, reviews_per_month, last_review</li> </ul>

## DATA UNDERSTANDING AND PREPARATION

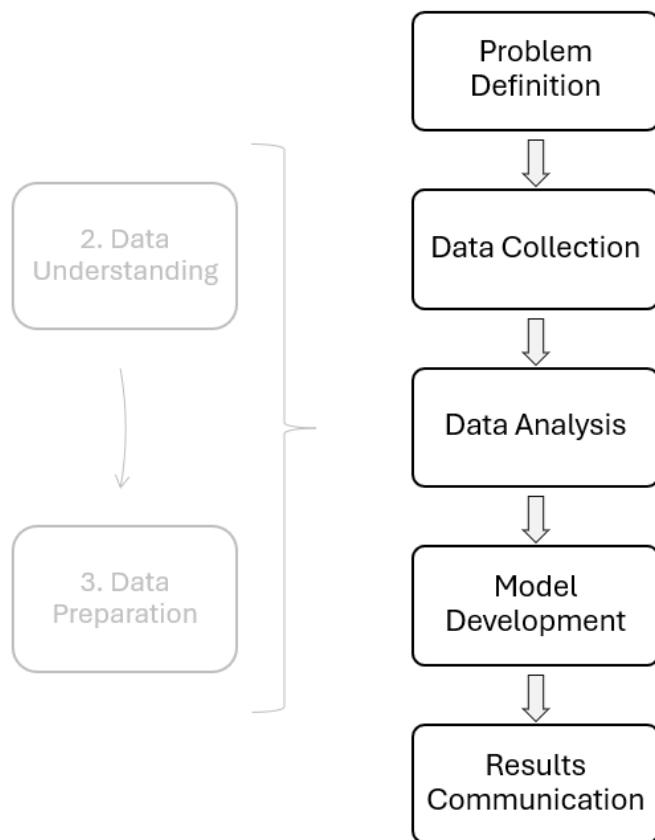


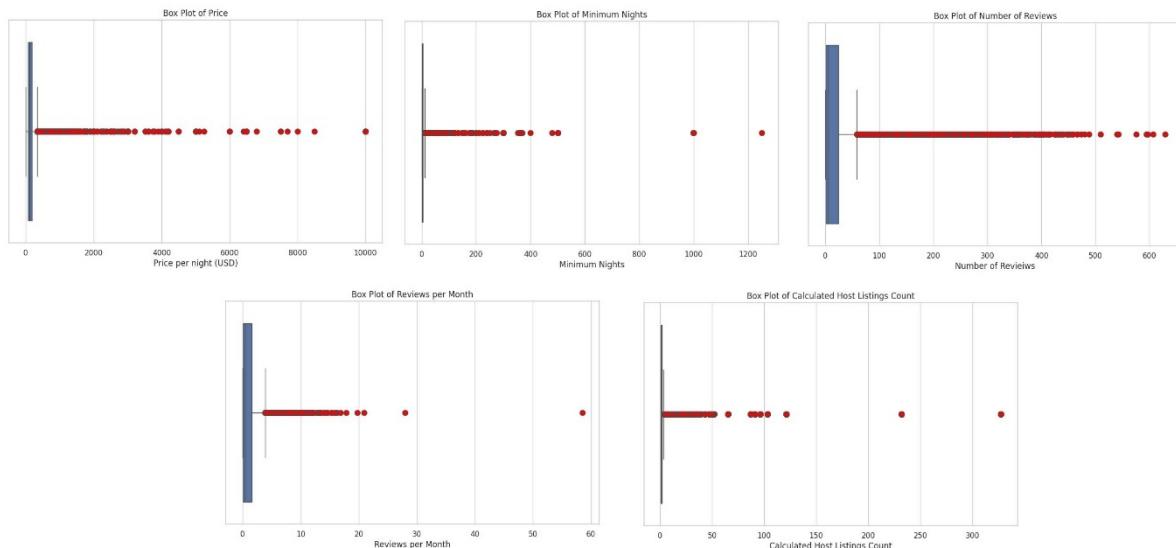
FIGURE 4 | EDA

EDA follows CRISP-DM (Figure 4) with data understanding identifying missing values (Table 2). (Mukhiya & Ahmed, 2020; Oluleye, 2023). Data preparation set missing `reviews_per_month` to 0 and imputed `last_review` with earliest date and adding a `has_review` flag. Non-sensitive Personally Identifiable Information (PII) `host_name` and `name` were removed (IBM, N.D.).

**TABLE 2 |** Missing data

<b>Count</b>	48,895 entries	
<b>Columns</b>	16 columns	10 numeric 6 categorical
<b>Missing data</b>	4 columns	last_review 10052 reviews_per_month 10052 host_name 21 name 16

Statistical analysis indicated non-normal distribution, confirmed by skew, kurtosis, and boxplots (Figure 5) (Mukhiya & Ahmed, 2020). Anderson-Darling confirmed highly significant deviation across all features at a 1% significance level (Berenson et al., 2019; Bobbitt, 2019). Because outliers shift the mean, Table 3 shows mode and median are more reliable measures of central tendency (Holmes et al., 2022). Outliers were removed using Interquartile Range (IQR), reducing the dataset to 29,408 entries (Bruce et al., 2020).



**FIGURE 5 |** Outliers (red) in box plots (price, minimum\_nights, number\_of\_reviews, reviews\_per\_month, and calculated\_host\_listings\_count)

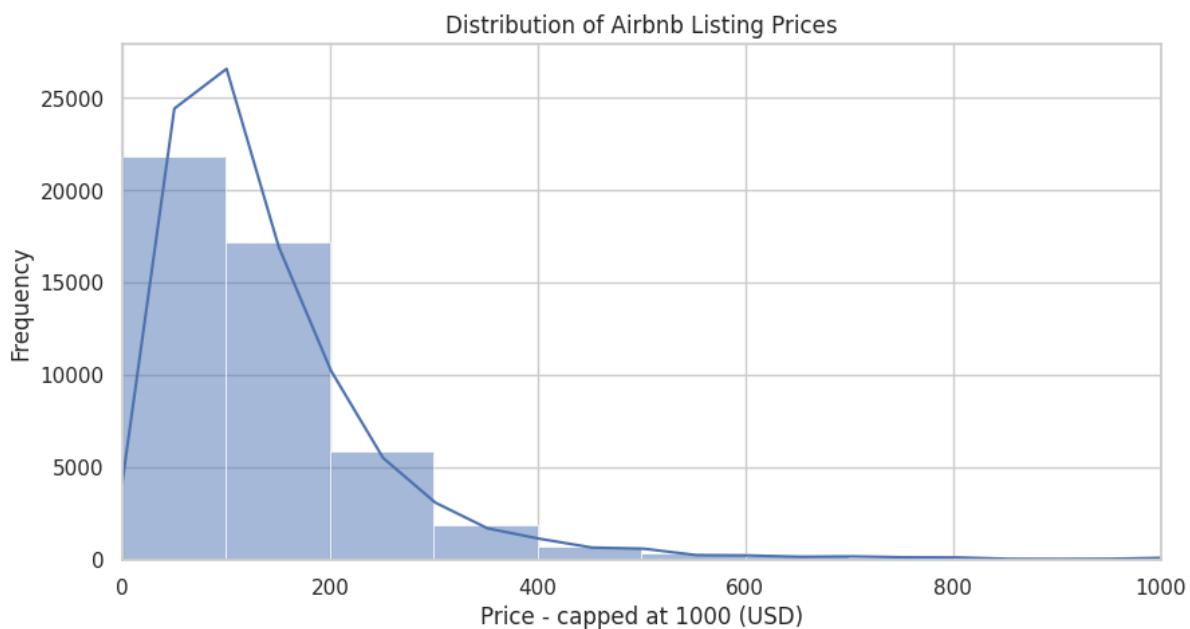
**TABLE 3 | Price**

Mode	\$100.00
Median	\$106.00
Mean	\$152.72
Min	\$0
Max	\$10,000
Upper band of price outliers	\$334.00
Number of outliers	2,972

## DATA LEARNINGS

### Price

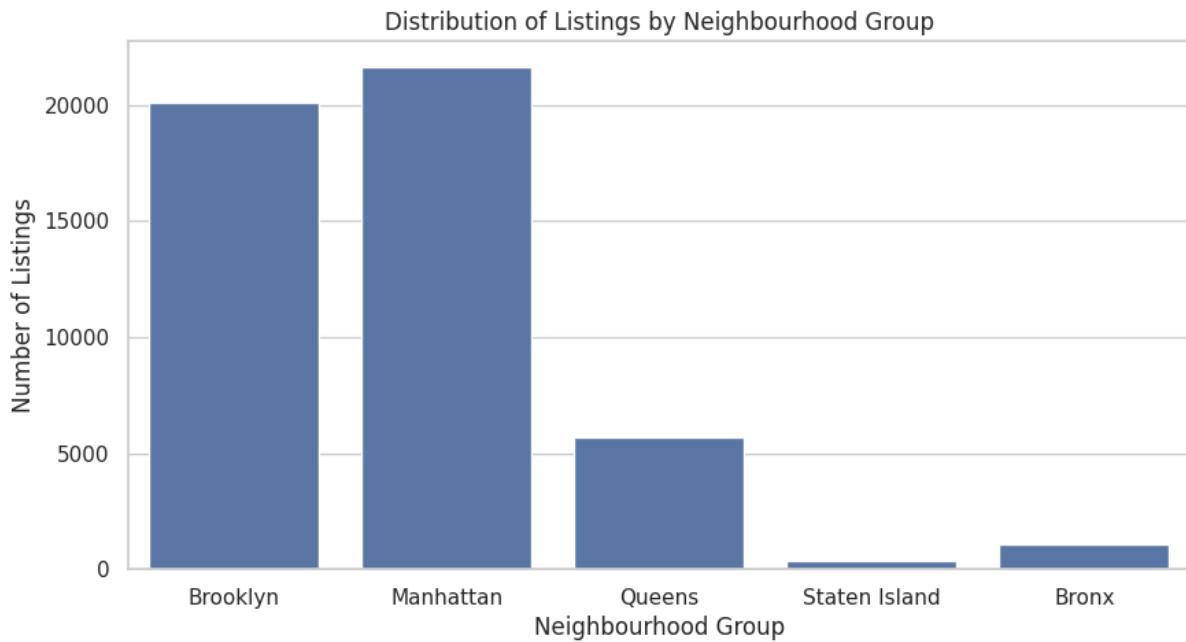
Figure 6 shows a modal peak at \$100 (Oluleye, 2023). The long right tail starts at \$334, with 2972 outliers reaching up to \$10,000. This indicates price variability from luxury listings.



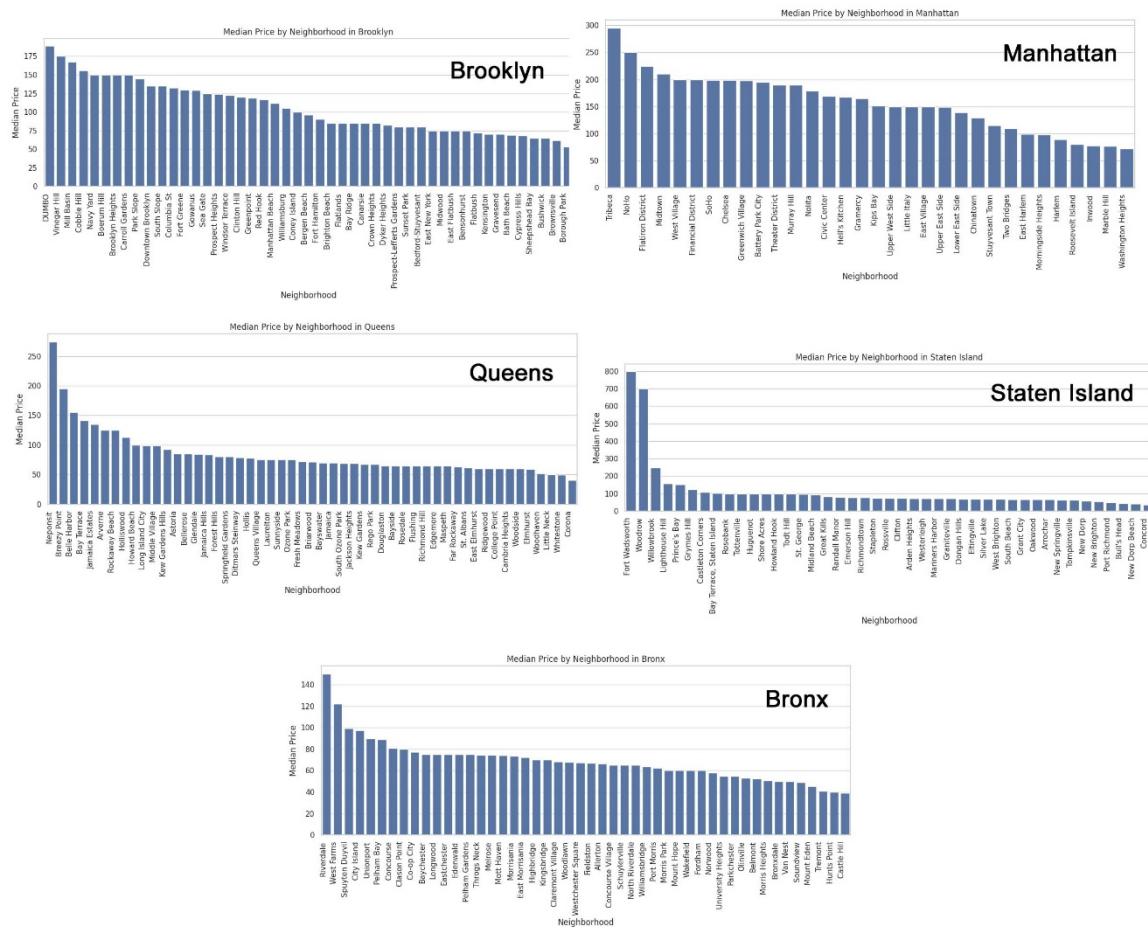
**FIGURE 6 | Price**

## Location

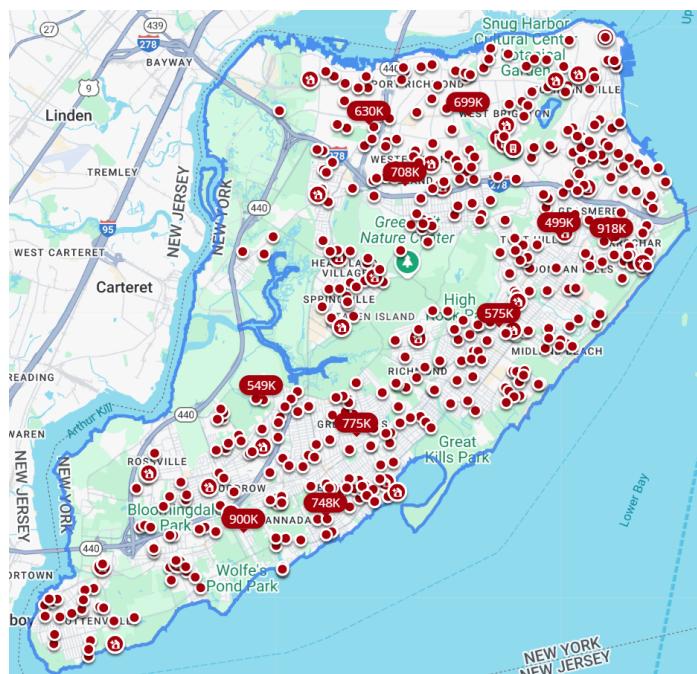
Most listings are in Manhattan and Brooklyn (Figure 7), with high density and pricing indicating strong demand in these popular boroughs (Figure 3). Median neighbourhood prices vary (Figure 8), with Staten Island's Fort Wadsworth and Woodrow reflecting higher Zillow (N.D.) property values (Figure 9).



**FIGURE 7 |** Neighbourhood Group



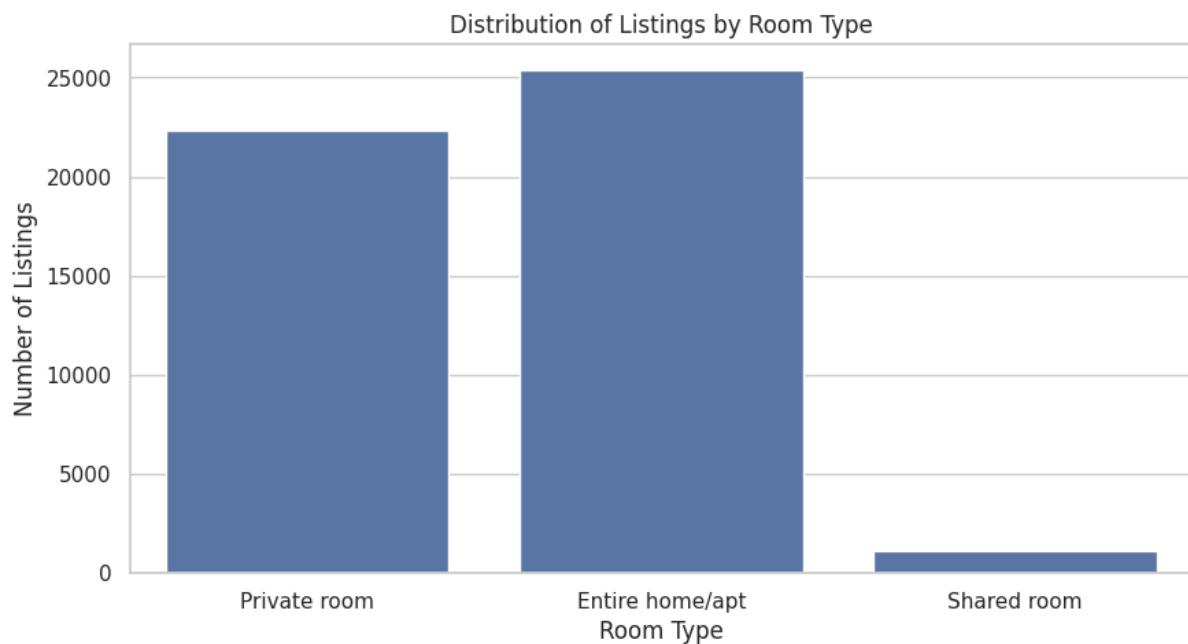
**FIGURE 8 | Neighbourhood median prices**



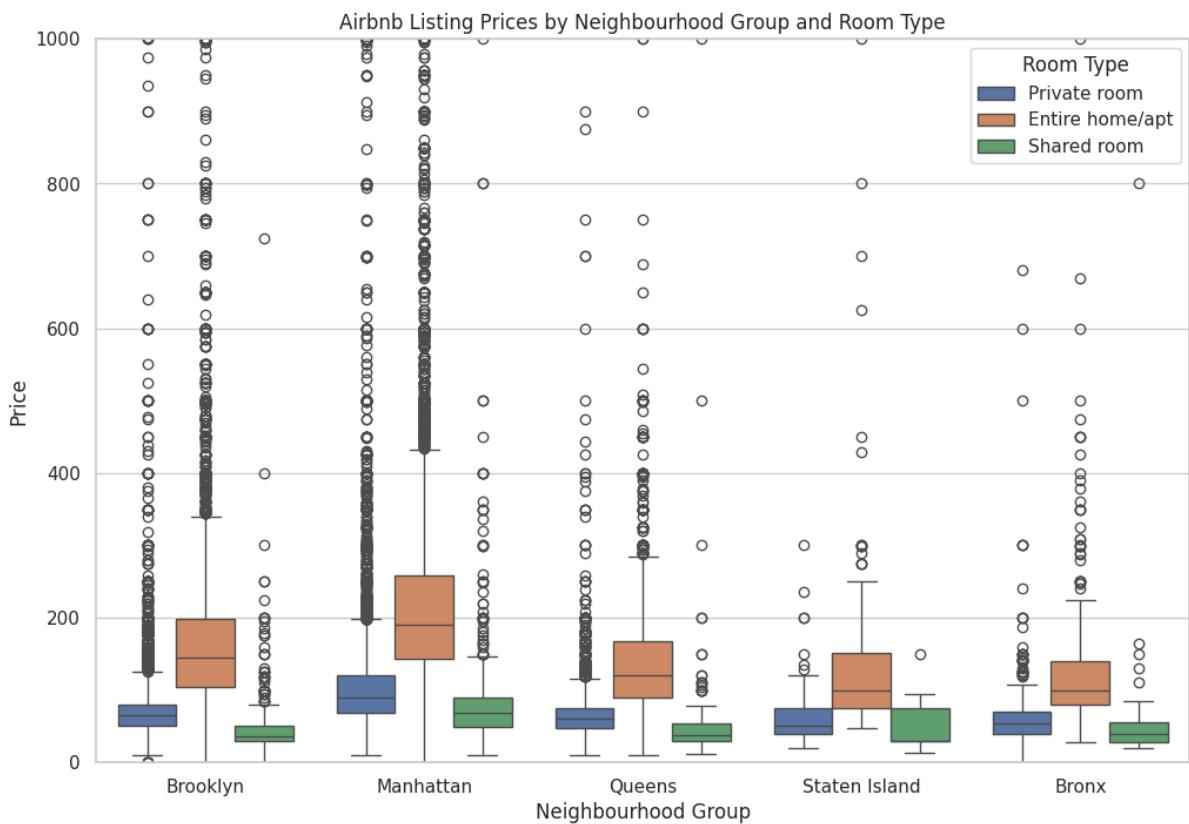
**FIGURE 9 | Staten Island property prices**

## Room Type

Most listings are entire homes, followed by private rooms (Figure 10). Manhattan's entire homes fetch the highest prices (Figure 11), highlighting demand for private, exclusive, central accommodation.



**FIGURE 10 | Room type**



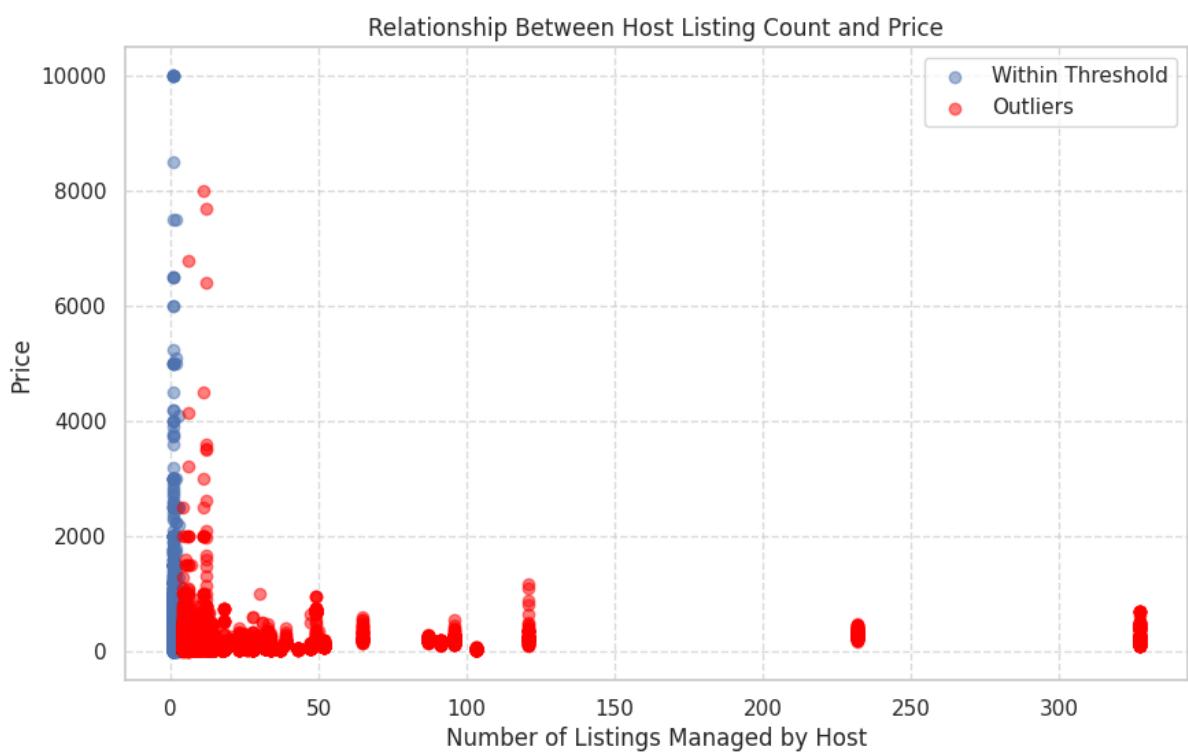
**FIGURE 11 |** Price per neighbourhood group by room type

## Host Characteristics

Of 37,457 hosts, 5,154 professionals manage multiple listings (Figure 12) (Abrate et al., 2022). Removing outliers above 3.5 listings reduces hosts with full-year availability (Figure 13, Table 4). Regardless of outliers, correlation between price and number of listings per host is weak.



**FIGURE 12 |** Single versus multiple listings (left), outliers removed (right)



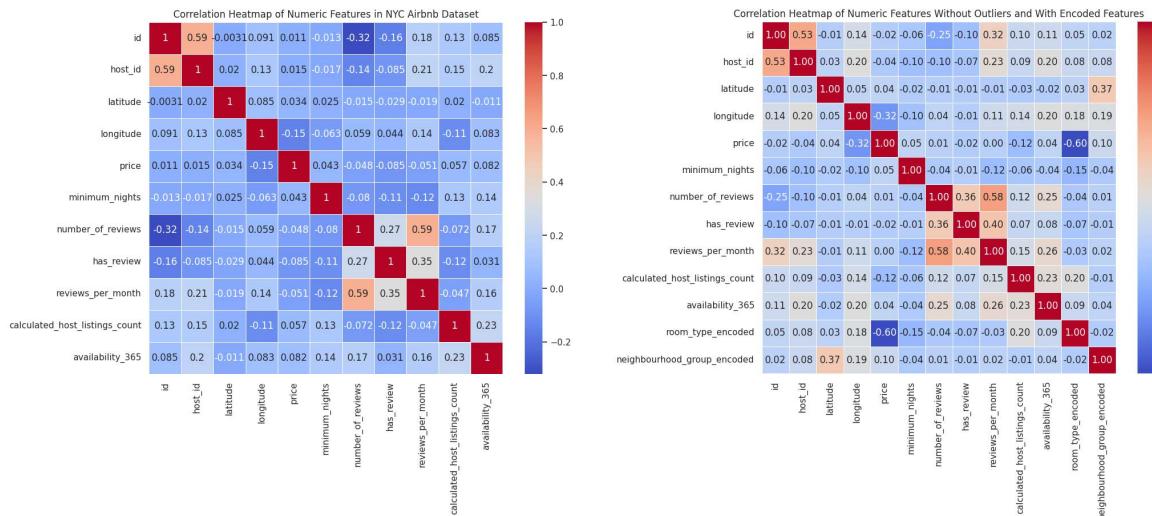
**FIGURE 13 |** Price per host listing count

**TABLE 4 |** Host characteristics

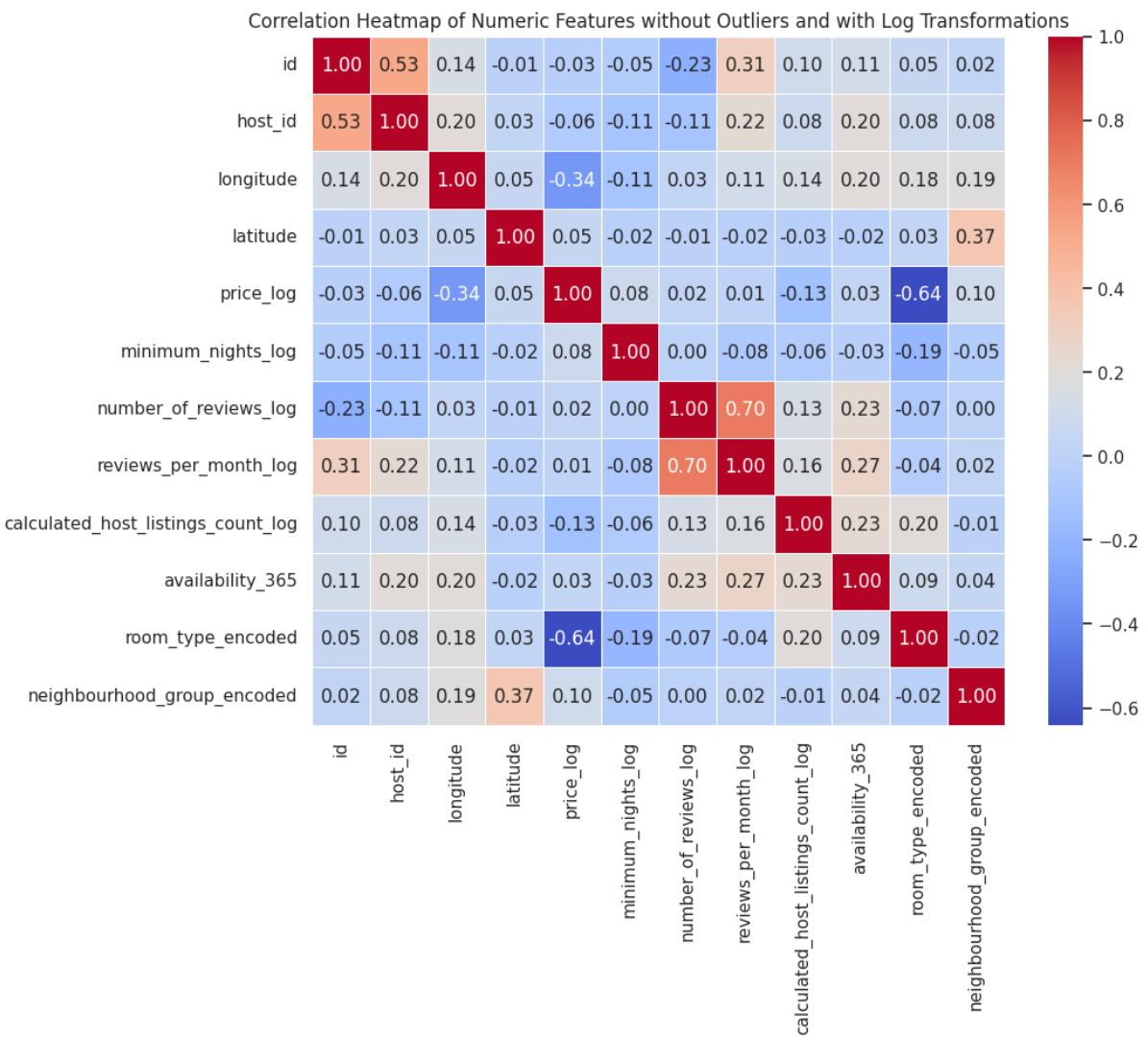
Metric	Before Outlier Removal	After Outlier Removal
Total Listings	48895.00	41814.00
Unique Hosts	37457.00	36583.00
Hosts with 1 Listing	32303.00	32303.00
Hosts with More Than 1 Listing	5154.00	4280.00
Average Reviews Per Host	1.03	1.02
Hosts with Full Year Availability	894.00	741.00

## Correlation

Figure 14 indicates weak correlations, with a notable negative correlation (-0.60), between `room_type_encoded` and `price`. Log transformation (Figure 13) had minimal impact on correlation strength.



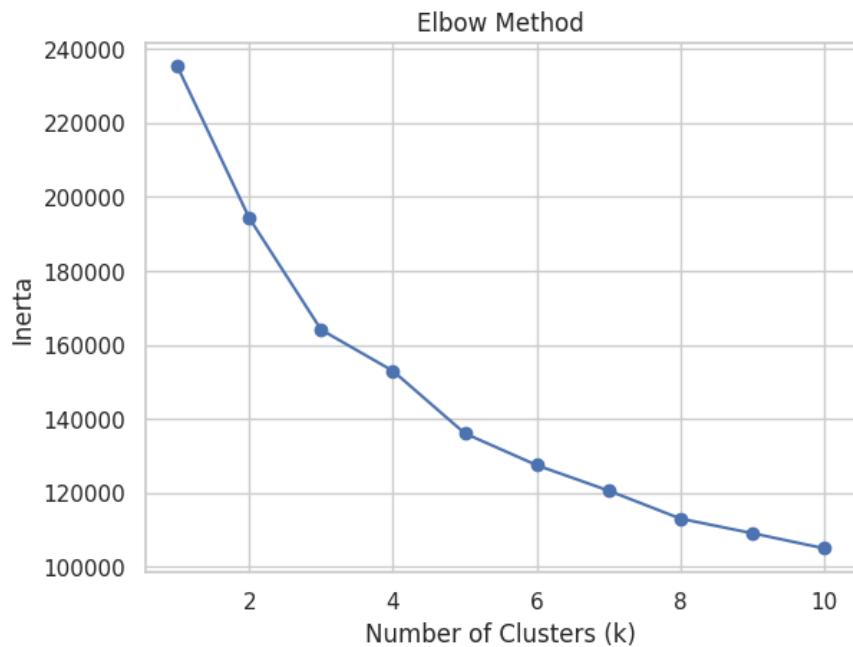
**FIGURE 14 |** Linear correlation (left); without outliers and with encoded categoricals (right)



**FIGURE 15 |** Log-linear correlation

## Clustering

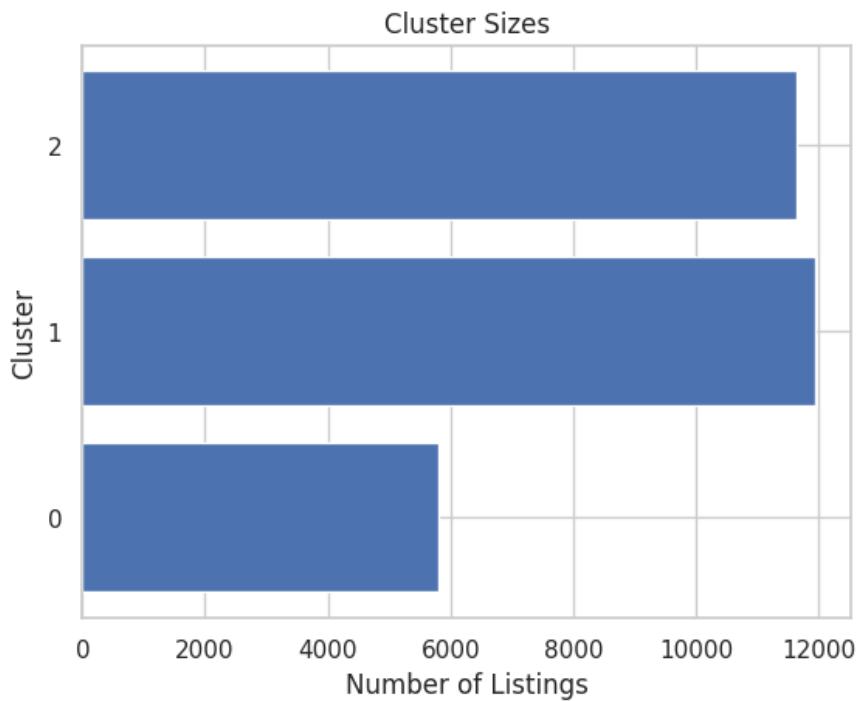
Three optimal  $k$ -means clusters were identified (Figure 16), with numerical encoding enhancing analysis (Table 5). Cluster 0, half the size, shows higher activity with the most reviews and availability (Figures 17, 18). Cluster 1 features higher pricing concentrated in Manhattan and Brooklyn, while Cluster 2 represents budget-friendly options spread throughout NYC (Figures 18, 19).



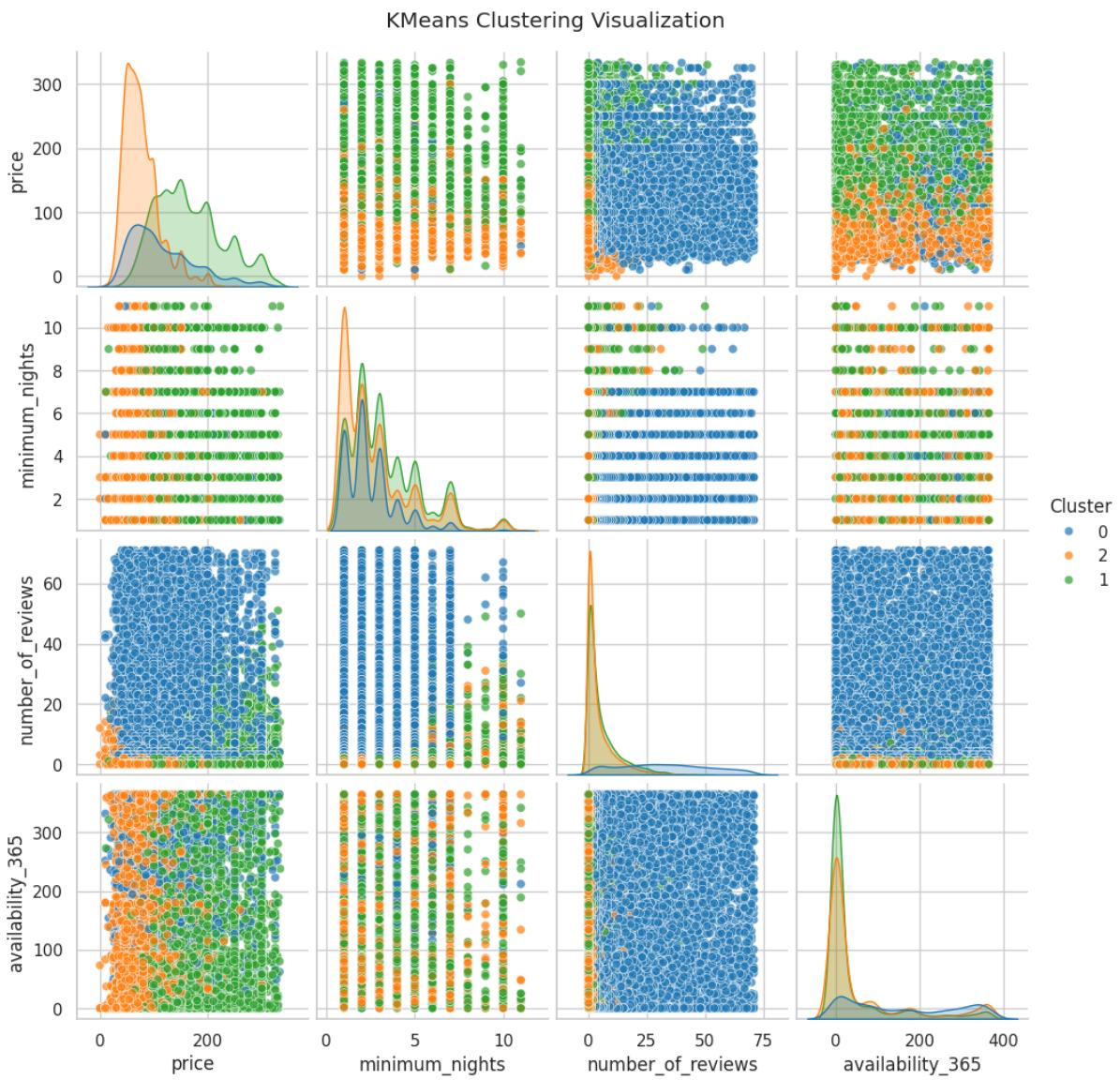
**FIGURE 16 |** Elbow method

**TABLE 5 |** Encoding

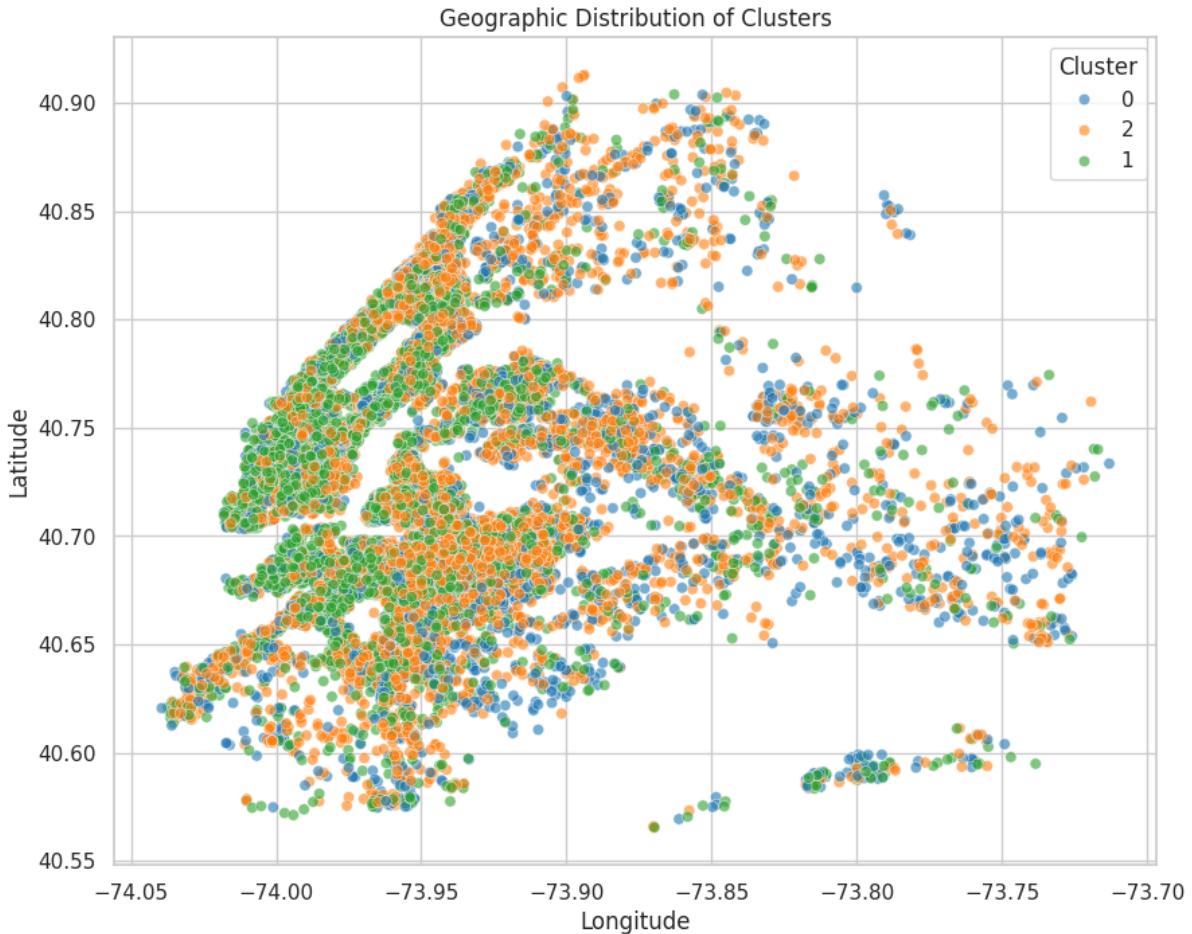
room_type	Entire: 0, Private: 1, Shared: 2
neighbourhood_group	Bronx: 0, Brooklyn: 1, Manhattan: 2, Queens: 3, Staten Island: 4



**FIGURE 17 |** Cluster sizes



**FIGURE 18 |** Clustering



**FIGURE 19 |** Geographic clustering

## INSIGHTS

Revenue optimisation requires cluster-specific strategies (Tables 6, 7). Cluster 1's exclusivity commands highest prices, but Cluster 0's moderate pricing and higher availability outpaces revenue at the same booking rate. Half of Cluster 0 and 1 bookings have reviews, while Cluster 2 has one-fifth, suggesting potential to improve relationships. Dynamic pricing benefits all hosts, particularly Cluster 1 during peaks, with Clusters 0 and 2 benefiting seasonally (Abrate et al., 2022). NYC's Local Law 18 and market saturation pose challenges as expansion markets outpace core markets (Airbnb, N.D.b; Voltes-Dorta, 2024). Opportunities exist through data integration, like census, and advanced ML for improved targeting and revenue strategies (Zhu et al., 2020; Akalın & Alptekin, 2024).

**TABLE 6 | Cluster mean analysis**

Cluster (Listings)	Price	Minimum Nights	Number of Reviews	Has Review	Reviews per Month	Calculated Host Listings Count	Availability 365	Room Type Encoded	Neighbourhood Group Encoded
0 Moderately priced, available (5807)	£117.37	2.5	33	1	1.73	1.5	156.2	0.53 Split entire home and private room	1.66 Majority Brooklyn, some Manhattan
1 High-priced homes, limited availability (11,934)	£166.34	3.3	6	0.77	0.33	1.11	40.5	0.02 Entire home	1.65 Majority Brooklyn, some Manhattan
2 Budget-friendly rooms, moderate availability (11,658)	£77.08	2.7	4	0.7	0.28	1.3	59.6	1.04 Mostly private room / few shared	1.57 Less Manhattan, Brooklyn, Bronx

**TABLE 7 | Example revenue**

	Price x Listings x Availability	Booking Rate (%)	ANNUAL REVENUE (USD millions)
Cluster 0	\$106.46m	50%	\$53.23m
Cluster 1	\$80.40m	50%	\$40.20m
Cluster 2	\$53.56m	50%	\$26.78m

## CONCLUSION

This analysis aimed to optimise Airbnb's NYC revenue through EDA and cluster-specific insights. Tailored pricing strategies—exclusivity for luxury, enhanced booking rates for moderate, improved engagement for budget, and dynamic pricing for all—can boost profitability. Challenges like Local Law 18 and market saturation underscore the need for continuous data-driven adaptation to continue revenue growth in the home-sharing market (Voltes-Dorta, 2024).

## REFERENCES

Abrate, G., Sainaghi, R. & Mauri, A.G. (2022) Dynamic pricing in Airbnb: Individual versus professional hosts, *Journal of Business Research* 141: 191–199. DOI: <https://doi.org/10.1016/J.JBUSRES.2021.12.012>.

Airbnb (N.D.a) *About Us*. Available from: <https://news.airbnb.com/about-us/> [Accessed 25 November 2024].

Airbnb (N.D.b) *Airbnb Investor: Quarterly Results*. Available from: <https://investors.airbnb.com/financials/default.aspx#quarterly> [Accessed 25 November 2024].

Akalın, Ö. & Alptekin, G.I. (2024) Enhancing Airbnb Price Predictions with Location-Based Data: A Case Study of Istanbul 207–212. DOI: <https://doi.org/10.15439/2024F7603>.

Berenson, M., Levine, D., Szabat, K. & Stephan, D. (2019) *Basic Business Statistics, Global Edition*. 14th ed. Harlow: Pearson International Content.

Bobbitt, Z. (2019) *How to Conduct an Anderson-Darling Test in R*. Available from: <https://www.statology.org/anderson-darling-test-r/> [Accessed 27 November 2024].

Bruce, P., Bruce, A. & Gedeck, P. (2020) *Practical Statistics for Data Scientists*. 2nd ed. O'Reilly Media. Available from: <vbk://9781492072898> [Accessed 24 May 2024].

Ding, K., Niu, Y. & Choo, W.C. (2023) The evolution of Airbnb research: A systematic literature review using structural topic modeling, *Helion* 9(6). DOI: <https://doi.org/10.1016/j.heliyon.2023.e17090>.

Gunter, U., Önder, I. & Zekan, B. (2020) Modeling Airbnb demand to New York City while employing spatial panel data at the listing level, *Tourism Management* 77: 104000. DOI: <https://doi.org/10.1016/J.TOURMAN.2019.104000>.

Holmes, A., Illowsky, B. & Dean, S. (2022) *Introductory Business Statistics*. OpenStax. Available from: <https://openstax.org/details/books/introductory-business-statistics> [Accessed 25 April 2024].

IBM (N.D.) *What is Personally Identifiable Information (PII)*? Available from: <https://www.ibm.com/topics/pii> [Accessed 27 November 2024].

Kaggle (2021) *New York City Airbnb Open Data*, Kaggle. Available from: <https://www.kaggle.com/code/whyalwaysme/ab-nyc-2019> [Accessed 26 November 2024].

Mukhiya, S.K. & Ahmed, U. (2020) *Hands-On Exploratory Data Analysis with Python*. Packt Publishing.

Niakšu, O. (2015) CRISP Data Mining Methodology Extension for Medical Domain, *Baltic J Modern Computing* 3(2): 92–109. Available from: [https://www.researchgate.net/publication/277775478\\_CRISP\\_Data\\_Mining\\_Methodology\\_Extension\\_for\\_Medical\\_Domain](https://www.researchgate.net/publication/277775478_CRISP_Data_Mining_Methodology_Extension_for_Medical_Domain) [Accessed 30 October 2023].

Oluleye, A. (2023) *Exploratory Data Analysis with Python Cookbook*. Pockt Publishing. Available from: <https://learning.oreilly.com/library/view/exploratory-data-analysis/9781803231105/> [Accessed 12 November 2024].

Taylor, B. & Almeida, A. (2022) *Revenue Growth, Stock Analysis*. Available from: <https://stockanalysis.com/term/revenue-growth/> [Accessed 26 November 2024].

Voltes-Dorta, A. (2024) Price competition among short-term Airbnb listings in New York City following Local Law 18, *Current Issues in Tourism* [Preprint]. DOI: <https://doi.org/10.1080/13683500.2024.2354521>.

Zhu, A., Li, R. & Xie, Z. (2020) Machine Learning Prediction of New York Airbnb Prices, *Proceedings - 2020 3rd International Conference on Artificial Intelligence for Industries, AI4I 2020* 1–5. DOI: <https://doi.org/10.1109/AI4I49448.2020.00007>.

Zillow (N.D.) *Staten Island Property Prices*. Available from: <https://www.zillow.com/staten-island-new-york-ny/> [Accessed 28 November 2024].

## APPENDIX



Code: <https://github.com/mariaingold/AirbnbNYC>

**kaggle** Dataset: [AB\\_NYC\\_2019.csv](#)



Colab:

<https://colab.research.google.com/github/mariaingold/AirbnbNYC/blob/main/AirbnbNYC.ipynb>

### MACHINE LEARNING: AIRBNB ANALYSIS

#### PROJECT

UoE Machine Learning

Assignment Due Date: 2 December, 2024

#### TEAM 2 AUTHORS

- Ahmed Husain
- Dinh (Danty) Khoi
- Maria Ingold
- Murthy Kanuri

#### PROJECT DESCRIPTION

##### Development Team Project

Determine interesting business analytic substantive question answered by dataset that is useful for the Airbnb executive board.

##### Business Analytic Question

Agreed by: Ahmed, Danty, Maria, Murthy

How can Airbnb NYC optimise revenue growth through tailored pricing recommendations based on location, room type, and host characteristics?

**Source Dataset: Airbnb NYC 2019**

#### Code

#### CONTRIBUTIONS

ACTION	WHO	DESCRIPTION
Setup	Maria	Google Colab, GitHub (connection, collaboration, Readme), structure, initial setup, basic data exploration.
EDA	Danty	Exploratory Data Analysis (EDA)
	Ahmed	Exploratory Data Analysis
	Maria	Basic EDA as part of setup, comment detail, outlier boxplots, added neighbourhood, added elbow method
	Murthy	Mode and Range
Machine Learning	Danty	Clustering

ACTION	WHO	DESCRIPTION
	Ahmed	Clustering
Review	Ahmed	EDA and clustering
	Danty	EDA and clustering
	Maria	Merged initial notebook with Danty's and Ahmed's EDA and clustering, Murthy's mode and range analysis, removed duplications, consolidated variable names, fixed warnings, added structure and additional comments.
	Murthy	EDA and clustering
Analysis	Ahmed	EDA and clustering
	Danty	EDA and clustering. Comment detail. Input into report draft.
	Maria	EDA and clustering. Detailed comments to provide analysis of results. Further research into outliers especially for pricing. Refined missing data handling. Added more on cluster mean analysis.
	Murthy	EDA and clustering. Input into report draft.

## SETUP

### Import Libraries

In [1]:

```
import pandas as pd          # data manipulation
import numpy as np           # numerical computation / linear algebra
import seaborn as sns         # enhanced visuals and statistics
import geopandas as gpd       # geospatial data
import matplotlib.lines as mlines   # customised legend lines / markers
import folium                 # interactive maps
from matplotlib import pyplot as plt # plots, charts and figures
from scipy import stats        # Anderson-Darling test for normality
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler, LabelEncoder, MinMaxScaler
from sklearn.metrics import silhouette_score
from shapely.geometry import Point
```

### Import Data

*Manually upload AB\_NYC\_2019.csv into Session Storage in Google Colab*

### Read Data

In [2]:

```
nyc_airbnb_df = pd.read_csv("AB_NYC_2019.csv")
```

## EXPLORATORY DATA ANALYSIS

### Overview

#### Column overview

The NYC Airbnb dataset has 48,895 entries and 16 columns. Key columns include:

- Host Information: host\_id and host\_name.
- Location Information: neighbourhood\_group, neighbourhood, latitude, longitude.
- Listing Characteristics: room\_type, price, minimum\_nights, number\_of\_reviews, last\_review, and reviews\_per\_month.
- Listing Metrics: calculated\_host\_listings\_count (number of listings by the host) and availability\_365 (days available per year).

Some columns, such as name, host\_name, and last\_review, contain missing values. Specifically:

- name and host\_name have a few missing entries.
- last\_review and reviews\_per\_month have more substantial missing values, likely due to properties with no reviews.

### EDA Plan for This Dataset

1. Data Understanding: Structure and contents: columns, rows, column names and types
2. Data Cleaning: Identify and handle missing values.
3. Descriptive Analysis: Descriptive statistical analysis like skew and kurtosis for numeric features.
4. Univariate Analysis: Examine distributions for numerical columns (e.g., price, minimum\_nights).
5. Bivariate Analysis: Examine two variables.

## 6. Multivariate Analysis: Clustering and more.

### Steps to Address Question

1. Analyze Price Distribution: Examine price to identify outliers and distribution shape.
2. Price by Neighbourhood Group: Analyze how prices vary by neighbourhood\_group.
3. Price by Room Type: Examine price trends across different room types.
4. Price vs Host Characteristics. Calculated Host Listings Count: Investigate if hosts with multiple listings set different price levels.
5. Interactions: Neighbourhood Group & Room Type: Analyze combined influence on prices.
6. Cluster analysis.

### DATA UNDERSTANDING

Get a sense of the data structure and contents by summarising the data.

#### Number of columns and rows

In [3]:

```
print("Number of columns: ", len(nyc_airbnb_df.columns))
print("Number of rows: ", len(nyc_airbnb_df))
```

Number of columns: 16

Number of rows: 48895

#### Column names, non-nulls and types

In [4]:

```
nyc_airbnb_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48895 entries, 0 to 48894
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   id               48895 non-null   int64  
 1   name              48879 non-null   object  
 2   host_id            48895 non-null   int64  
 3   host_name           48874 non-null   object  
 4   neighbourhood_group 48895 non-null   object  
 5   neighbourhood        48895 non-null   object  
 6   latitude             48895 non-null   float64 
 7   longitude            48895 non-null   float64 
 8   room_type             48895 non-null   object  
 9   price                48895 non-null   int64  
 10  minimum_nights       48895 non-null   int64  
 11  number_of_reviews     48895 non-null   int64  
 12  last_review           38843 non-null   object  
 13  reviews_per_month      38843 non-null   float64 
 14  calculated_host_listings_count 48895 non-null   int64  
 15  availability_365       48895 non-null   int64  
dtypes: float64(3), int64(7), object(6)
memory usage: 6.0+ MB
```

#### 5-Figure Summary per Column

- Mean (mean)
- Standard Deviation (std)
- Minimum and Maximum (min, max)
- Median (50%)
- Quartile 1 (25%)
- Quartile 3 (75%)

In [5]:

```
nyc_airbnb_df.describe()
```

Out[5]:

	<b>id</b>	<b>host_id</b>	<b>latitude</b>	<b>longitude</b>	<b>price</b>	<b>minimum_nights</b>	<b>number_of_reviews</b>	<b>reviews_per_month</b>	<b>calculated_host_listings_count</b>	<b>availability_365</b>
<b>co</b>	4.889	4.889	48895	48895	48895	48895	48895.0	38843.0	48895.000000	48895.000000
<b>un</b>	500e+	500e+	.0000	.0000	.0000	.0000	00000	00000		
<b>t</b>	04	04	00	00	00	00				
<b>me</b>	1.901	6.762	40.72	-	152.7	7.029	23.2744	1.37322	7.143982	112.78
<b>an</b>	714e+	001e+	8949	73.95	20687	962	66	1		1327
	07	07		2170						

	<b>id</b>	<b>host_id</b>	<b>latitude</b>	<b>longitude</b>	<b>price</b>	<b>minimum_nights</b>	<b>number_of_reviews</b>	<b>reviews_per_month</b>	<b>calculated_host_listings_count</b>	<b>availability_365</b>
<b>std</b>	1.098 311e+07	7.861 097e+07	0.054 530	0.046 157	240.1 54170	20.51 0550	44.5505 82	1.68044 2	32.952519	131.62 2289
<b>min</b>	2.539 000e+03	2.438 000e+03	40.49 9790	-74.24 4420	0.000 000	1.000 000	0.000000 0	0.01000 0	1.000000	0.0000 00
<b>25%</b>	9.471 945e+06	7.822 033e+06	40.69 0100	-73.98 3070	69.00 0000	1.000 000	1.000000 0	0.19000 0	1.000000	0.0000 00
<b>50%</b>	1.967 728e+07	3.079 382e+07	40.72 3070	-73.95 5680	106.0 00000	3.000 000	5.000000 0	0.72000 0	1.000000	45.000 000
<b>75%</b>	2.915 218e+07	1.074 344e+08	40.76 3115	-73.93 6275	175.0 00000	5.000 000	24.0000 00	2.02000 0	2.000000	227.00 0000
<b>max</b>	3.648 724e+07	2.743 213e+08	40.91 3060	-73.71 2990	10000 .0000	1250. 00000	629.000 000	58.5000 00	327.000000	365.00 0000

#### Mode

This is the first mode for each column.

In [6]:

```
# Calculate the mode for each column
mode_values = nyc_airbnb_df.mode().iloc[0]
```

```
# Display the mode values
print(mode_values)
```

<b>id</b>	2539
<b>name</b>	Hillside Hotel
<b>host_id</b>	219517861.0
<b>host_name</b>	Michael
<b>neighbourhood_group</b>	Manhattan
<b>neighbourhood</b>	Williamsburg
<b>latitude</b>	40.71813
<b>longitude</b>	-73.95677
<b>room_type</b>	Entire home/apt
<b>price</b>	100.0
<b>minimum_nights</b>	1.0
<b>number_of_reviews</b>	0.0
<b>last_review</b>	2019-06-23
<b>reviews_per_month</b>	0.02
<b>calculated_host_listings_count</b>	1.0
<b>availability_365</b>	0.0

Name: 0, dtype: object

#### Range

In [7]:

```
# Calculate the range for each numerical column
range_values = nyc_airbnb_df.select_dtypes(include='number').apply(lambda x: x.max() - x.min())
```

```
# Convert the range values to avoid exponential notation
range_values_formatted = range_values.apply(lambda x: f"{x:.0f}" if x >= 1 else f"{x:.6f}")
```

```
# Display the range values
print("Range for each numerical column:")
print(range_values_formatted)
```

Range for each numerical column:

<b>id</b>	36484706
<b>host_id</b>	274318875
<b>latitude</b>	0.413270

```

longitude          0.531430
price              10000
minimum_nights     1249
number_of_reviews   629
reviews_per_month   58
calculated_host_listings_count 326
availability_365    365
dtype: object
Review Data (head and tail)
Reviewing both head and tail shows a more complete picture.
Head
Head shows mostly complete data.
In [8] :

nyc_airbnb_df.head()

Out[8] :

      h
      o  ho  neig  nei  la  lo  ro  p  min  numb  la
      s  st  hbou  ghb  ti  ng  om  r  imu  er_o  st
      d  e  t  _n  rhoo  our  tu  it  _t  i  m_n  f_re  ews_
      e  _am  d_gr  hoo  de  ud  yp  c  igh  view  _r  per_
      i  e  oup  d  e  e  ts  s  ie  h  mont  t_listi  calcla
      d                    w  nt  ngs_cou  avai
                           d                  labi
                           d                  lity
                           d                  _365

      Cle
      an
      &
      qui
      2 et  2
      0 5 apt 7 Jo Broo Ken sin 40 73 - Pr
      3 hom 8 hn klyn gto gto 47 72 at 4 1 9 20
      9 e   7
      by
      the
      par
      k

      Sky
      lit 2 Je Manh Mid 40 73 - En
      1 5 Mid 8 nn atta tow .7 .9 re 2 20
      9 tow 4 if n n 53 83 ho 2 1 45 19
      5 n Cas 5 er n n 62 77 me 5 05
      tle
      /a
      pt 21

      THE
      VIL
      LAG
      E
      3 OF 4 El Manh Har 40 73 - Pr
      2 6 HAR 6 is Manh Har .8 .9 at 1
      4 LEM 3 ab atta lem 09 41 e 5 3 0 Na
      7 ... 2 et n n 02 90 ro 0 NaN
      .NE h
      .om
      W
      YOR
      K !

      Coz
      y
      Ent
      3 ire 4 Li Cli 40 73 - En
      8 Flo 8 sa Broo nto 40 73 - ti
      3 or 6 Ro Broo n .6 .9 re 8 20
      1 of 9 xa klyn Hil 85 59 ho 9 19
      Bro 9 nn Hil 14 76 me 9 07
      wns e l /a 270 4.64 1 194
      ton
      e
      pt 05

      4 5 Ent 7 La Manh Eas 40 73 - En
      0 ire 1 ur atta t .7 .9 ti 8 20
      Apt 1 a n re 0 10 9 18 0.10 1 0
      -
```



		h	o	ho	neig	nei	a	l	lo	ro	p	min	numb	la	revi	calcula	ava
i	nam	s	st	hbou	ghb	i	ng	om	r	imu	er_o	st	ews_	ted_hos	ila		
d	e	t	_n	rhoo	our	t	it	_t	i	m_n	f_re	r	per_	t_listi	bil		
		am	d_gr	hoo	ud	ud	yp	c	igh	view	ev	mont	ngs_cou	ity			
		i	e	oup	d	e	e	e	ts	s	ie	h	nt		36		
		d			d	e					w				5		

Sun  
ny  
Stu  
3 dio 2 Il 4 En  
4 at 3 ga 0 - ti  
His 9 r Manh . 73 re 1  
8 tor 2 & atta Har 8 . 9 ho 1 10 0 Na NaN 1 27  
ica 9 Ay n lem 1 48 me 5 N N  
l 5 se 4 48 /a  
3 Nei 1 7 67 pt  
1 ghb 2 5 5  
orh  
ood

43r  
d  
St.  
3 Tim 3 4  
4 e 0 Hel 0 - Sh  
Squ 9 8 Ta Manh l's 7 73 ar 5 1 0 Na NaN 6 2  
8 are 8 5 z atta Kit 5 . 9 ed 5 1 ro 0 N N  
- coz 7 n che 7 91 ro om  
3 Y 9 sin 5 1  
gle  
bed

Tre  
ndy  
dup  
lex  
3 in 6 4  
4 the 8 Ch 0 - Pr  
ver 1 . Hel .  
8 y 1 ri Manh l's 7 73 iv  
8 hea 9 st atta Kit 6 . 9 at 9 7 0 Na NaN 1 23  
9 op 8 n che 4 89 e 0 ro N N  
4 rt 8 he 4 33 om  
of 1 n 0  
5 Hel 4 4  
l's  
Kit  
che  
n

#### DATA CLEANING

##### Identify missing values

In [10]:

```
# Identify missing values in the dataset
missing_values = nyc_airbnb_df.isnull().sum()
missing_values = missing_values[missing_values > 0].sort_values(ascending=False)
missing_values
```

Out[10]:

0

last_review	10052
reviews_per_month	10052
host_name	21
name	16

```

dtype: int64
reviews_per_month
Missing reviews likely means none have been made, so set to 0.
In [11]:


# Set missing reviews per month to 0
nyc_airbnb_df.loc[nyc_airbnb_df['reviews_per_month'].isnull(), 'reviews_per_month'] = 0

last_review

- Missing last review likely means there is none. Set date to be before 2008 where it can't exist, but keeps datetime type.
- Don't want to lose that it has no review so add has_review column.


In [12]:


# Fill NaNs with a dummy date outside the range to keep datetime dtype
nyc_airbnb_df['last_review'] = nyc_airbnb_df['last_review'].fillna(pd.Timestamp.min)

# Add column for feature engineering with 0 for no review, or 1 for has review.
nyc_airbnb_df['has_review'] = nyc_airbnb_df['last_review'].apply(lambda x: 0 if x == pd.Timestamp.min else 1)

# Reorder columns to place 'has_review' after 'last_review'
cols = list(nyc_airbnb_df.columns)
last_review_index = cols.index('last_review')
cols.insert(last_review_index + 1, cols.pop(cols.index('has_review')))
nyc_airbnb_df = nyc_airbnb_df.reindex(columns=cols)

host_name and name

- host_name is personally identifiable information (PII) and should be removed.
- name (of listing) is not being used for this analysis, so has been removed.


In [13]:


# drop host_name and name
nyc_airbnb_df = nyc_airbnb_df.drop(['host_name', 'name'], axis=1, errors='ignore')

Verify missing values after handling
In [14]:


missing_values_after_cleaning = nyc_airbnb_df.isnull().sum()
missing_values_after_cleaning

Out[14]:

```

	0
<b>id</b>	0
<b>host_id</b>	0
<b>neighbourhood_group</b>	0
<b>neighbourhood</b>	0
<b>latitude</b>	0
<b>longitude</b>	0
<b>room_type</b>	0
<b>price</b>	0
<b>minimum_nights</b>	0
<b>number_of_reviews</b>	0
<b>last_review</b>	0
<b>has_review</b>	0

```

          0
reviews_per_month      0
calculated_host_listings_count   0
availability_365        0

dtype: int64
In [15]:
nyc_airbnb_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48895 entries, 0 to 48894
Data columns (total 15 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               48895 non-null   int64  
 1   host_id          48895 non-null   int64  
 2   neighbourhood_group 48895 non-null   object  
 3   neighbourhood     48895 non-null   object  
 4   latitude          48895 non-null   float64 
 5   longitude         48895 non-null   float64 
 6   room_type         48895 non-null   object  
 7   price             48895 non-null   int64  
 8   minimum_nights    48895 non-null   int64  
 9   number_of_reviews 48895 non-null   int64  
 10  last_review       48895 non-null   object  
 11  has_review        48895 non-null   int64  
 12  reviews_per_month 48895 non-null   float64 
 13  calculated_host_listings_count 48895 non-null   int64  
 14  availability_365   48895 non-null   int64  
dtypes: float64(3), int64(8), object(4)
memory usage: 5.6+ MB
In [16]:
nyc_airbnb_df.head()

Out[16]:
   h
   o   neigh   nei   la   lo   ro   p   min   numb   ha   revi   calculat   avai
   s   bourh   ghb   ti   ng   om   r   imu   er_o   last   s_   ews_   ed_host_
   t   ood_g   our   tu   it   _t   i   m_n   f_re   _rev   re   per_   listings   labi
   d   ood_g   hoo   tu   ud   yp   c   igh   view   iew   vi   mont   _count   lity
   i   roup    d   de   e   e   e   ts   s   ew   h
   d

      Pr
 2  2   Ken  40  -  iv  1
 0  5  7  Brook  sin .6  .9  at  4  1  9  2018
 0  3  8  lyn  gto  47  72  e  9
 0  9  7   n  49  37  ro  om
      En
 2  2
 1  5  8  Manha  Mid  40  -  ti
 1  9  4  ttan   tow  .7  .9  re  2
 1  5  5   n   53  83  ho  2  1  45  2019
 1  5  5   n   62  77  me  5  -05-
 1  5  5   n   77  /a  21
 1  5  5   n   77  pt

      Pr
 2  3  4
 2  6  6  Manha  Har  40  -  iv  1  1677
 2  4  3  ttan   lem  .8  .9  at  5  -09-
 2  4  3  ttan   lem  09  .9  e  0  21
 2  7  2   n  41  41  ro  0  00:1
 2  7  2   n  90  90  om  0  2:43
 2  7  2   n  93  93  145
 2  7  2   n  93  93  2241
 2  7  2   n  93  93  93


```

```

    h
  o   neigh   nei   la   lo   ro   p   min   numb   ha   revi   calculat   avai
  i   s   neigh   nei   la   lo   ro   p   min   numb   ha   revi   calculat   avai
  d   t   bourh   ghb   our   ti   ng   om   r   imu   er_o   last   s   ews_   ed_host_
  _i   _d   ood_g   hoo   tu   tu   it   ud   yp   c   igh   view   iew   re   per_   listings
  _i   _d   roup   d   de   e   e   ts   s   s   s   s   ew   h   mont   count   lity
  _d

            En
  3   4           Cli   40   -   re   8   1   2019
  8   8   Brook   nto   .6   .9   ho   9   -07-
  3   6   lyn    Hil   85   59   me   0   05
  1   9           l   14   76   /a
                           pt

            En
  5   7           Eas   40   -   re   8   10   9   2018
  0   1   Manha   t   .7   .9   ho   0   -11-
  2   9   ttan    Har   98   43   me   19
  2   2           lem   51   99   /a
                           pt

```

## Numeric Variables

In [17]:

```
# Selecting only numeric columns for analysis
numeric_columns_nyc = nyc_airbnb_df.select_dtypes(include=['float64', 'int64']).columns
```

## DESCRIPTIVE ANALYTICS

### Skew and Kurtosis

Descriptive analysis on numerical columns to examine distributions and calculate skewness and kurtosis for each.

#### SKEWNESS

- Asymmetry of distribution of dataset.
- Positive skew: right skewed (long right tail)
- Negative skew: left skewed (long left tail)
- Zero skew: normal distribution

High Positive Skew (large outliers towards right)

- price
- minimum\_nights
- number\_of\_reviews
- reviews\_per\_month
- calculated\_host\_listings\_count

Positive Skew (some outliers towards right)

- host\_id
  - longitude
- Slight Positive Skew (closer to normal, some outliers to right)
- latitude
  - availability\_365

Very Slightly Negative (almost normal)

- id

#### KURTOSIS

- Peakedness and tailedness of dataset distribution
- Positive kurtosis: peaked with more tails (outliers)
- Negative kurtosis: flatter and lighter tails (concentrated at mean)
- Zero kurtosis: normal peakedness in distribution

Very High Postive Kurtosis (sharp peak, heavy tails --> outliers)

- price
- minimum\_nights
- reviews\_per\_month
- calculated\_host\_listings\_count

High Positive Kurtosis (sharp peak, heavey tails --> outliers)

- longitude

Slight Positive Kurtosis (closer to normal, but outliers)

- host\_id

- latitude
- Sightly Negative Kurtosis
- id
  - has\_review
  - availability\_365

These findings suggest the need for potential transformations or outlier handling in these variables.

In [18]:

```
# Calculation of skewness and kurtosis for numerical columns

# Calculate skewness and kurtosis
skew_kurtosis_nyc = nyc_airbnb_df[numeric_columns_nyc].agg(['skew', 'kurtosis']).transpose()

# Display skewness and kurtosis results
skew_kurtosis_nyc
```

Out[18]:

	skew	kurtosis
<b>id</b>	-0.090257	-1.227748
<b>host_id</b>	1.206214	0.169106
<b>latitude</b>	0.237167	0.148845
<b>longitude</b>	1.284210	5.021646
<b>price</b>	19.118939	585.672879
<b>minimum_nights</b>	21.827275	854.071662
<b>number_of_reviews</b>	3.690635	19.529788
<b>has_review</b>	-1.457094	0.123127
<b>reviews_per_month</b>	3.300723	43.531611
<b>calculated_host_listings_count</b>	7.933174	67.550888
<b>availability_365</b>	0.763408	-0.997534

#### UNIVARIATE ANALYSIS

Olueye (2023) for visualising one variable:

- summary table
- bar chart
- boxplot (good for finding outliers)
- histogram (good for finding outliers)
- pie chart
- violin plot

#### Categorical Analysis

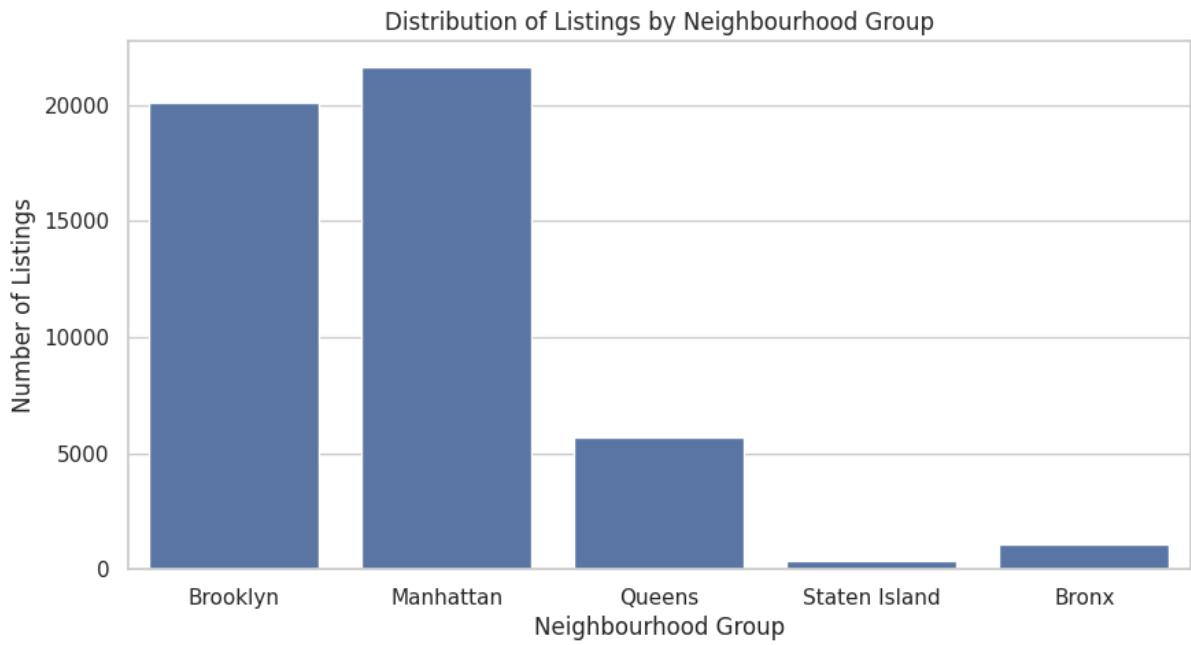
##### Neighbourhood Group

The majority of listings are concentrated in Manhattan and Brooklyn, with fewer in Queens, the Bronx, and Staten Island.

In [19]:

```
# Set a general style for the plots
sns.set(style="whitegrid")

# Plotting the distribution of listings by neighbourhood group
plt.figure(figsize=(10, 5))
sns.countplot(x='neighbourhood_group', data=nyc_airbnb_df)
plt.title("Distribution of Listings by Neighbourhood Group")
plt.xlabel("Neighbourhood Group")
plt.ylabel("Number of Listings")
plt.show()
```



#### **Unique Neighbourhoods in Neighbourhood Groups**

In [20]:

```
# Group by neighborhood_group and get unique neighborhoods
neighborhood_counts = nyc_airbnb_df.groupby('neighbourhood_group')['neighbourhood'].nunique()

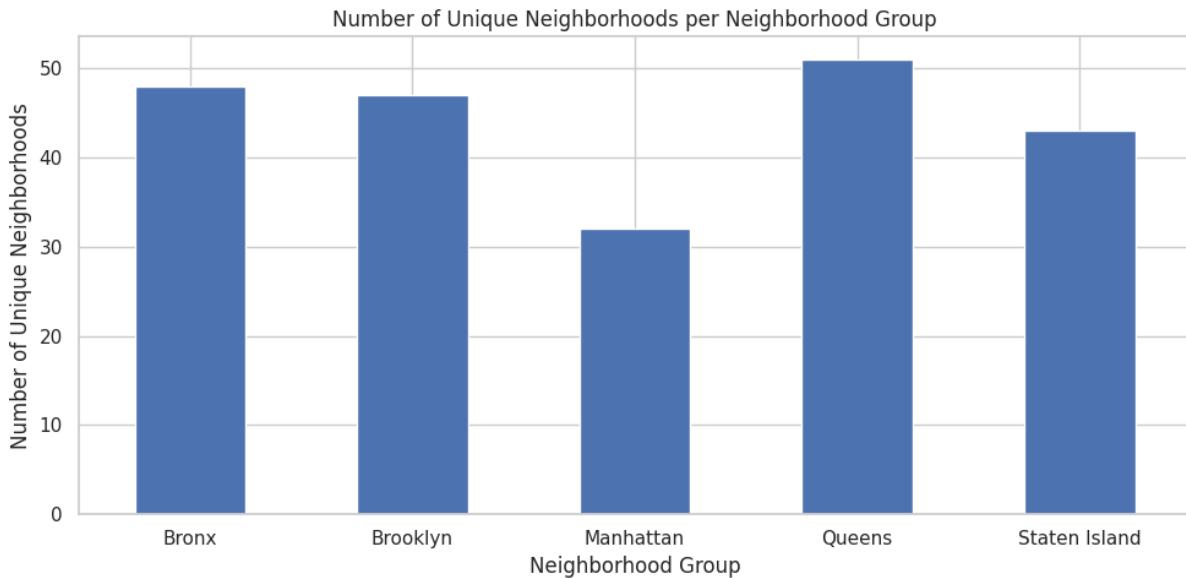
# Print the results
print("Unique Neighborhoods per Neighborhood Group:\n", neighborhood_counts)

# Visualization using a bar chart
plt.figure(figsize=(10, 5))
neighborhood_counts.plot(kind='bar')
plt.title('Number of Unique Neighborhoods per Neighborhood Group')
plt.xlabel('Neighborhood Group')
plt.ylabel('Number of Unique Neighborhoods')
plt.xticks(rotation=0) # Keep x-axis labels horizontal
plt.tight_layout()
plt.show()
```

Unique Neighborhoods per Neighborhood Group:

neighbourhood_group	
Bronx	48
Brooklyn	47
Manhattan	32
Queens	51
Staten Island	43

Name: neighbourhood, dtype: int64

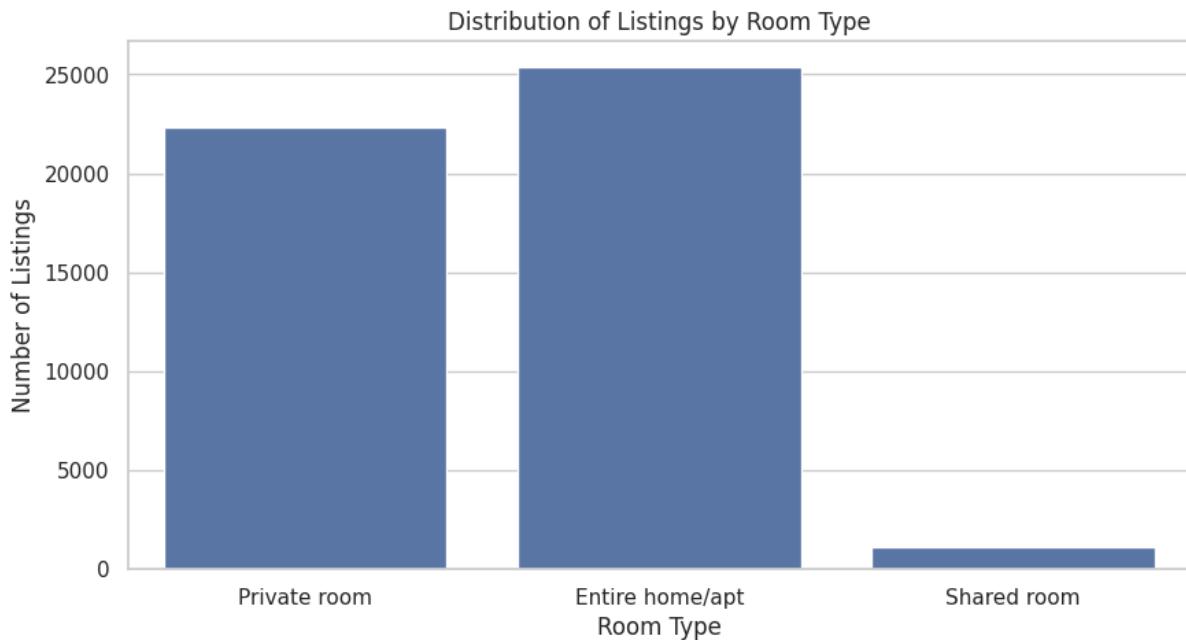


#### **Room Type**

Most listings are for "Entire home/apt," followed by "Private room," with fewer listings for "Shared room."

In [21]:

```
# Plotting the distribution of listings by room type
plt.figure(figsize=(10, 5))
sns.countplot(x='room_type', data=nyc_airbnb_df)
plt.title("Distribution of Listings by Room Type")
plt.xlabel("Room Type")
plt.ylabel("Number of Listings")
plt.show()
```

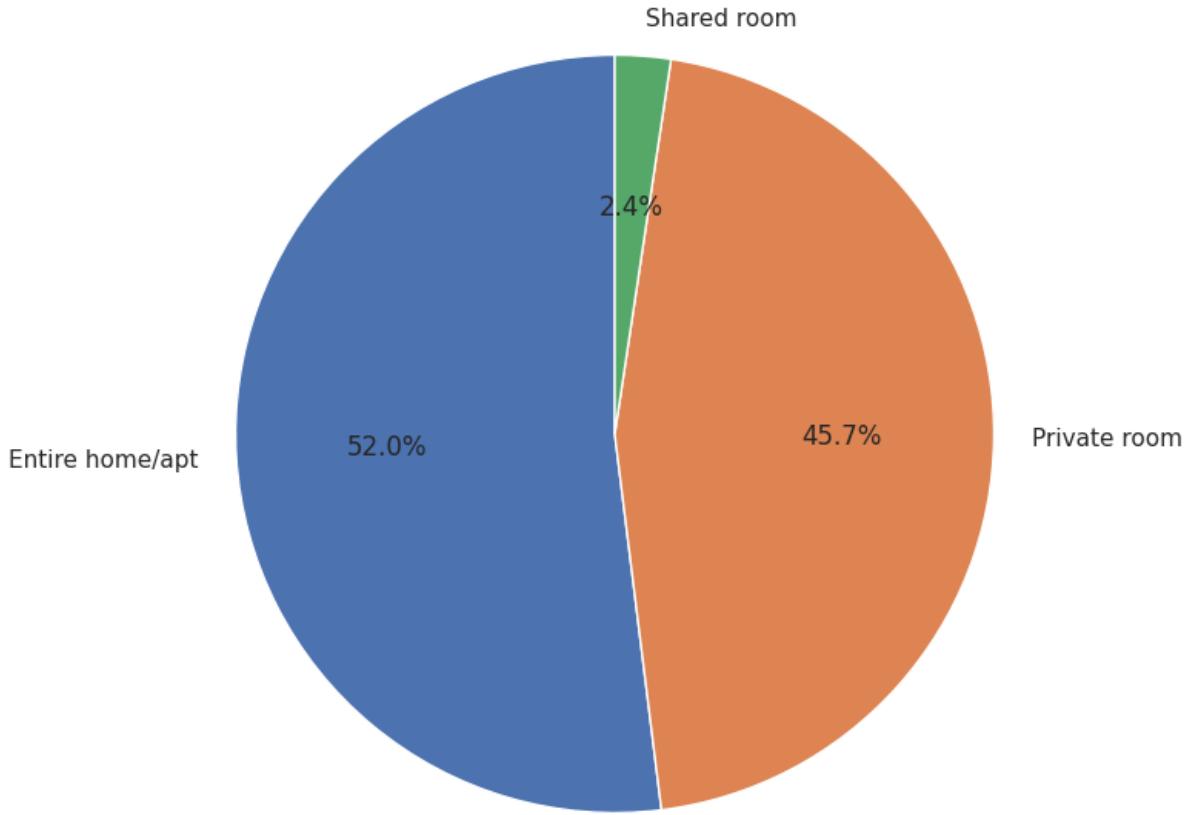


In [22]:

```
# Calculate the counts for each room type
room_type_counts = nyc_airbnb_df['room_type'].value_counts()

# Create the pie chart
plt.figure(figsize=(8, 8)) # Adjust figure size as needed
plt.pie(room_type_counts, labels=room_type_counts.index, autopct='%1.1f%%', startangle=90)
plt.title("Distribution of Listings by Room Type")
plt.show()
```

Distribution of Listings by Room Type



**Numerical Analysis**

**Calculate host listings count manually**

In [23]:

```
# Group by 'host_id' and count the number of listings for each host
host_listings_count = nyc_airbnb_df.groupby('host_id')['id'].count()

# Create a new column to classify hosts as having a single listing or multiple listings
single_listing = host_listings_count[host_listings_count == 1]
multiple_listings = host_listings_count[host_listings_count > 1]

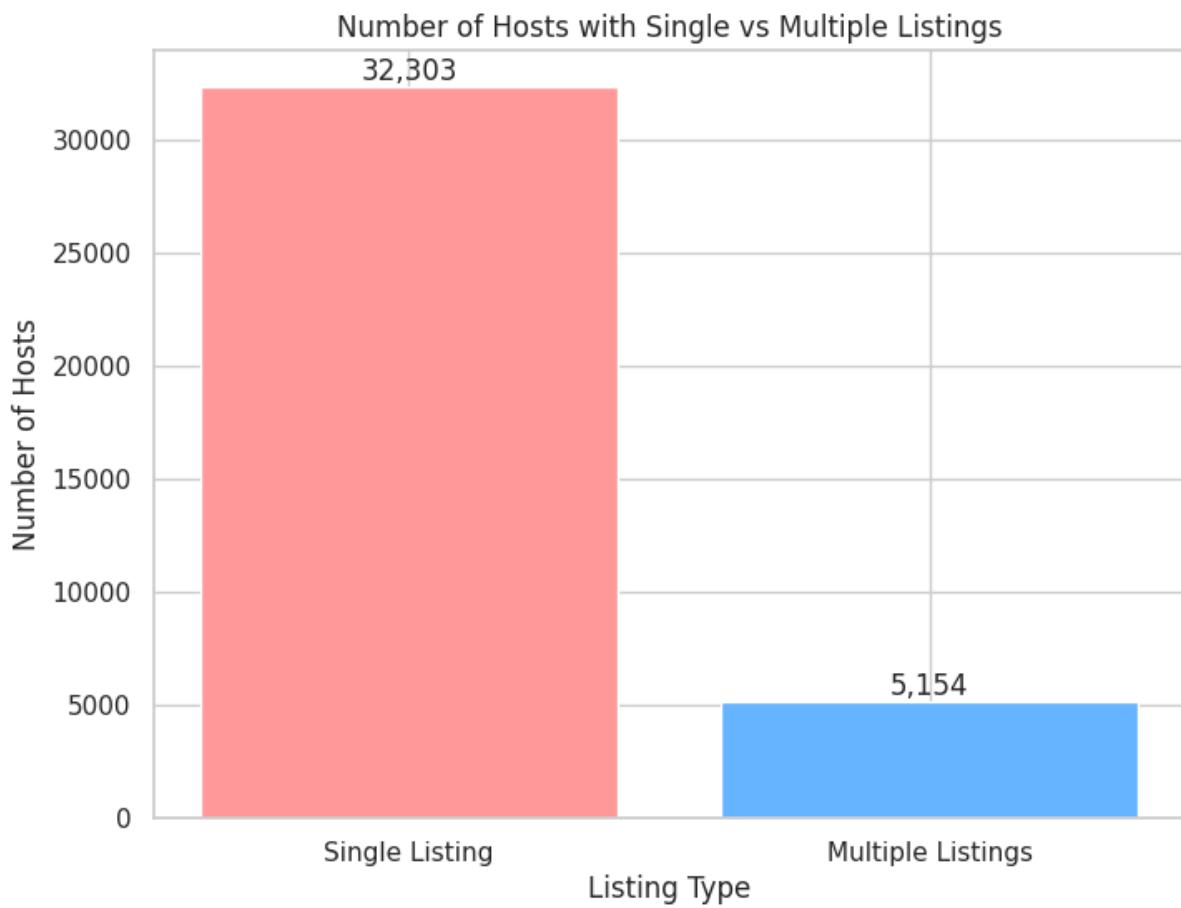
# Calculate counts
single_listing_count = len(single_listing)
multiple_listings_count = len(multiple_listings)

# Bar chart data
labels = ['Single Listing', 'Multiple Listings']
counts = [single_listing_count, multiple_listings_count]

# Create bar chart
plt.figure(figsize=(8, 6)) # Adjust figure size as needed
plt.bar(labels, counts, color=['#ff9999', '#66b3ff'])
plt.title('Number of Hosts with Single vs Multiple Listings')
plt.xlabel('Listing Type')
plt.ylabel('Number of Hosts')

# Add count labels on top of bars
for i, v in enumerate(counts):
    plt.text(i, v + 2, f'{v:,}', ha='center', va='bottom') # Adjust positioning as needed

plt.show()
```



In [24]:

```
# Group by 'host_id' and count the number of listings for each host
host_listings_count = nyc_airbnb_df.groupby('host_id')['id'].count()
```

```
print(host_listings_count.value_counts().sort_index())
```

id	Count
1	32303
2	3329
3	951
4	360
5	169
6	95
7	57
8	52
9	26
10	21
11	10
12	15
13	10
14	5
15	5
16	1
17	4
18	3
19	1
20	2
21	1
23	3
25	2
26	1
27	1
28	2
29	1
30	1
31	2
32	1

```

33      3
34      2
37      1
39      1
43      1
47      1
49      2
50      1
52      2
65      1
87      1
91      1
96      2
103     1
121     1
232     1
327     1
Name: count, dtype: int64
Calculate outliers for number of listings per host
There are so many at 1 that below or above 1 is an outlier.
In [25]:
host_listings_count.info()

<class 'pandas.core.series.Series'>
Index: 37457 entries, 2438 to 274321313
Series name: id
Non-Null Count Dtype
-----
37457 non-null int64
dtypes: int64(1)
memory usage: 585.3 KB
In [26]:
host_listings_count.describe()

Out[26]:
      id
count    37457.000000
mean     1.305363
std      2.760747
min     1.000000
25%    1.000000
50%    1.000000
75%    1.000000
max     327.000000
dtype: float64
In [27]:
# Calculate outlier thresholds using IQR
Q1 = host_listings_count.quantile(0.25)
Q3 = host_listings_count.quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Print outlier thresholds
print(f"Lower Outlier Threshold: {lower_bound}")
print(f"Upper Outlier Threshold: {upper_bound}")

Lower Outlier Threshold: 1.0
Upper Outlier Threshold: 1.0
Calculated Host Listings Count

```

- As a note, calculated\_host\_listings\_count shows that there are 327 listings with 1 host that has 327 listings denoted for each entry. So there is only one host with 327, but if count, then it appears 327 times. So, if need to count, can't use it like this.

In [28]:

```
print(nyc_airbnb_df['calculated_host_listings_count'].value_counts().sort_index())
```

calculated\_host\_listings\_count

1	32303
2	6658
3	2853
4	1440
5	845
6	570
7	399
8	416
9	234
10	210
11	110
12	180
13	130
14	70
15	75
16	16
17	68
18	54
19	19
20	40
21	21
23	69
25	50
26	26
27	27
28	56
29	29
30	30
31	62
32	32
33	99
34	68
37	37
39	39
43	43
47	47
49	98
50	50
52	104
65	65
87	87
91	91
96	192
103	103
121	121
232	232
327	327

Name: count, dtype: int64

**The number of listings that have hosts with this many listings.**

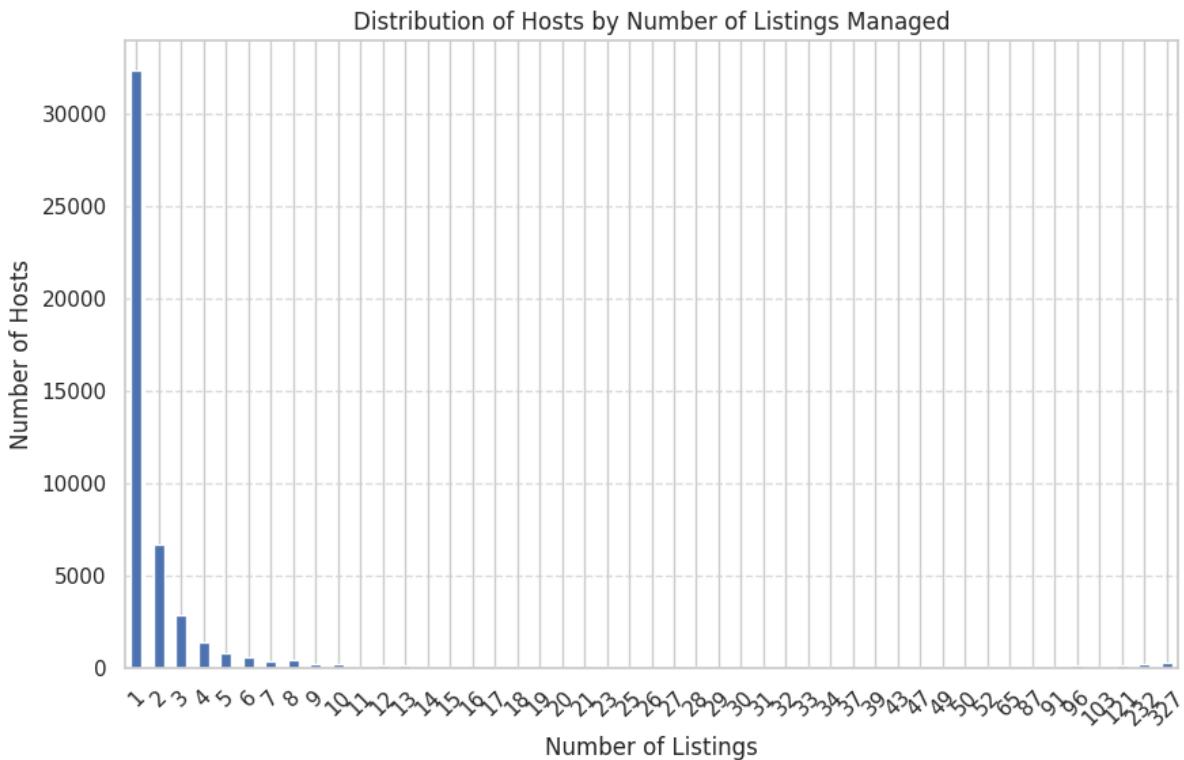
- 32,303 listings have 1 listing per host.
- 327 listings have 327 listed but this is all for 1 host.

This isn't a particularly useful visualisation.

In [29]:

```
# Distribution of hosts based on the number of listings they manage
host_listings_dist = nyc_airbnb_df['calculated_host_listings_count'].value_counts()

plt.figure(figsize=(10, 6))
host_listings_dist.sort_index().plot(kind='bar')
plt.title('Distribution of Hosts by Number of Listings Managed')
plt.xlabel('Number of Listings')
plt.ylabel('Number of Hosts')
plt.xticks(rotation=45)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.show()
```



In [30]:

```
calc_host_listings_count = nyc_airbnb_df['calculated_host_listings_count']
calc_host_listings_count.describe()
```

Out[30]:

	calculated_host_listings_count
<b>count</b>	48895.000000
<b>mean</b>	7.143982
<b>std</b>	32.952519
<b>min</b>	1.000000
<b>25%</b>	1.000000
<b>50%</b>	1.000000
<b>75%</b>	2.000000
<b>max</b>	327.000000
<b>dtype:</b>	float64

**Remove Outliers**

In [31]:

```
# Calculate quartiles
Q1 = nyc_airbnb_df['calculated_host_listings_count'].quantile(0.25)
Q3 = nyc_airbnb_df['calculated_host_listings_count'].quantile(0.75)
IQR = Q3 - Q1

# Calculate upper and lower bounds
upper_bound = Q3 + 1.5 * IQR
lower_bound = Q1 - 1.5 * IQR

print(f"Lower Outlier Threshold: {lower_bound}")
print(f"Upper Outlier Threshold: {upper_bound}")
```

```

# Filter data based on IQR bounds
filtered_data_iqr = nyc_airbnb_df[
    (nyc_airbnb_df['calculated_host_listings_count'] >= lower_bound) &
    (nyc_airbnb_df['calculated_host_listings_count'] <= upper_bound)
]

Lower Outlier Threshold: -0.5
Upper Outlier Threshold: 3.5
In [32]: 

def calculate_host_characteristics(df):
    """Calculates host characteristics for a given DataFrame."""
    unique_hosts = df['host_id'].nunique()
    multi_listing_hosts = df[df['calculated_host_listings_count'] > 1]['host_id'].nunique()
    hosts_with_one_listing = unique_hosts - multi_listing_hosts
    average_reviews_per_host = df.groupby('host_id')['reviews_per_month'].mean().mean()
    always_available_hosts = df[df['availability_365'] == 365]['host_id'].nunique()

    return {
        "Total Listings": len(df),
        "Unique Hosts": unique_hosts,
        "Hosts with 1 Listing": hosts_with_one_listing,
        "Hosts with More Than 1 Listing": multi_listing_hosts,
        "Average Reviews Per Host": average_reviews_per_host,
        "Hosts with Full Year Availability": always_available_hosts
    }

# Calculate characteristics for original data
original_characteristics = calculate_host_characteristics(nyc_airbnb_df)

# Calculate characteristics for filtered data
filtered_characteristics = calculate_host_characteristics(filtered_data_iqr)

# Set display options to prevent line wrapping
pd.set_option('display.width', None)
pd.set_option('display.max_colwidth', None)

# Create comparison table
comparison_table = pd.DataFrame({
    "Metric": list(original_characteristics.keys()),
    "Before Outlier Removal": list(original_characteristics.values()),
    "After Outlier Removal": list(filtered_characteristics.values())
})

# Format float columns to 2 decimal places
comparison_table_formatted = comparison_table.style.format({
    "Before Outlier Removal": "{:.2f}",
    "After Outlier Removal": "{:.2f}"
})

# Display the formatted table using display()
display(comparison_table_formatted)

```

Metric	Before Outlier Removal	After Outlier Removal
0 Total Listings	48895.00	41814.00
1 Unique Hosts	37457.00	36583.00
2 Hosts with 1 Listing	32303.00	32303.00
3 Hosts with More Than 1 Listing	5154.00	4280.00
4 Average Reviews Per Host	1.03	1.02
5 Hosts with Full Year Availability	894.00	741.00

#### Price by Listings Count (STEP 4)

In [33]:

```

# Set the upper outlier threshold
upper_threshold = upper_bound # Adjust this value as needed

# Scatter plot with outlier threshold

```

```

plt.figure(figsize=(10, 6))

# Plot data points within the threshold
plt.scatter(
    nyc_airbnb_df[nyc_airbnb_df['calculated_host_listings_count'] <=
upper_threshold]['calculated_host_listings_count'],
    nyc_airbnb_df[nyc_airbnb_df['calculated_host_listings_count'] <=
upper_threshold]['price'],
    alpha=0.5,
    label='Within Threshold'
)

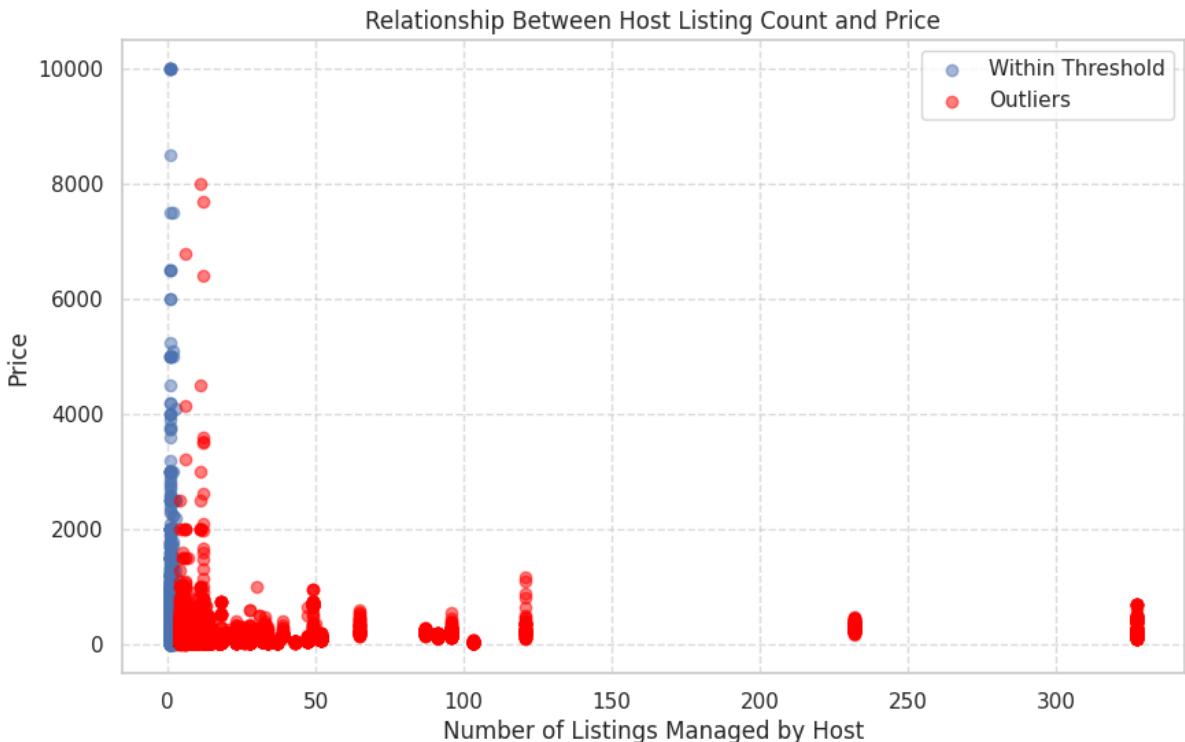
# Plot data points beyond the threshold (outliers)
plt.scatter(
    nyc_airbnb_df[nyc_airbnb_df['calculated_host_listings_count'] >
upper_threshold]['calculated_host_listings_count'],
    nyc_airbnb_df[nyc_airbnb_df['calculated_host_listings_count'] > upper_threshold]['price'],
    alpha=0.5,
    color='red', # Highlight outliers in red
    label='Outliers'
)

plt.title('Relationship Between Host Listing Count and Price')
plt.xlabel('Number of Listings Managed by Host')
plt.ylabel('Price')
plt.grid(linestyle='--', alpha=0.7)
plt.legend() # Add legend to distinguish data points
plt.show()

# Correlation between host listing count and price
correlation = nyc_airbnb_df[['calculated_host_listings_count', 'price']].corr().iloc[0, 1]
print("Correlation:", correlation)

# Correlation below the outlier threshold
filtered_data = nyc_airbnb_df[nyc_airbnb_df['calculated_host_listings_count'] <=
upper_threshold]
correlation_below_threshold = filtered_data[['calculated_host_listings_count',
'price']].corr().iloc[0, 1]
print("Correlation below threshold:", correlation_below_threshold)

```



```

Correlation: 0.057471688368068104
Correlation below threshold: -0.055767553095915426
In [34]:

```

```
# Group by 'host_id' and count the number of listings for each host
```

```

host_listings_count_filtered = filtered_data_iqr.groupby('host_id')['id'].count()

# Create a new column to classify hosts as having a single listing or multiple listings
single_listing_filtered = host_listings_count_filtered[host_listings_count_filtered == 1]
multiple_listings_filtered = host_listings_count_filtered[host_listings_count_filtered > 1]

# Calculate counts
single_listing_count_filtered = len(single_listing_filtered)
multiple_listings_count_filtered = len(multiple_listings_filtered)

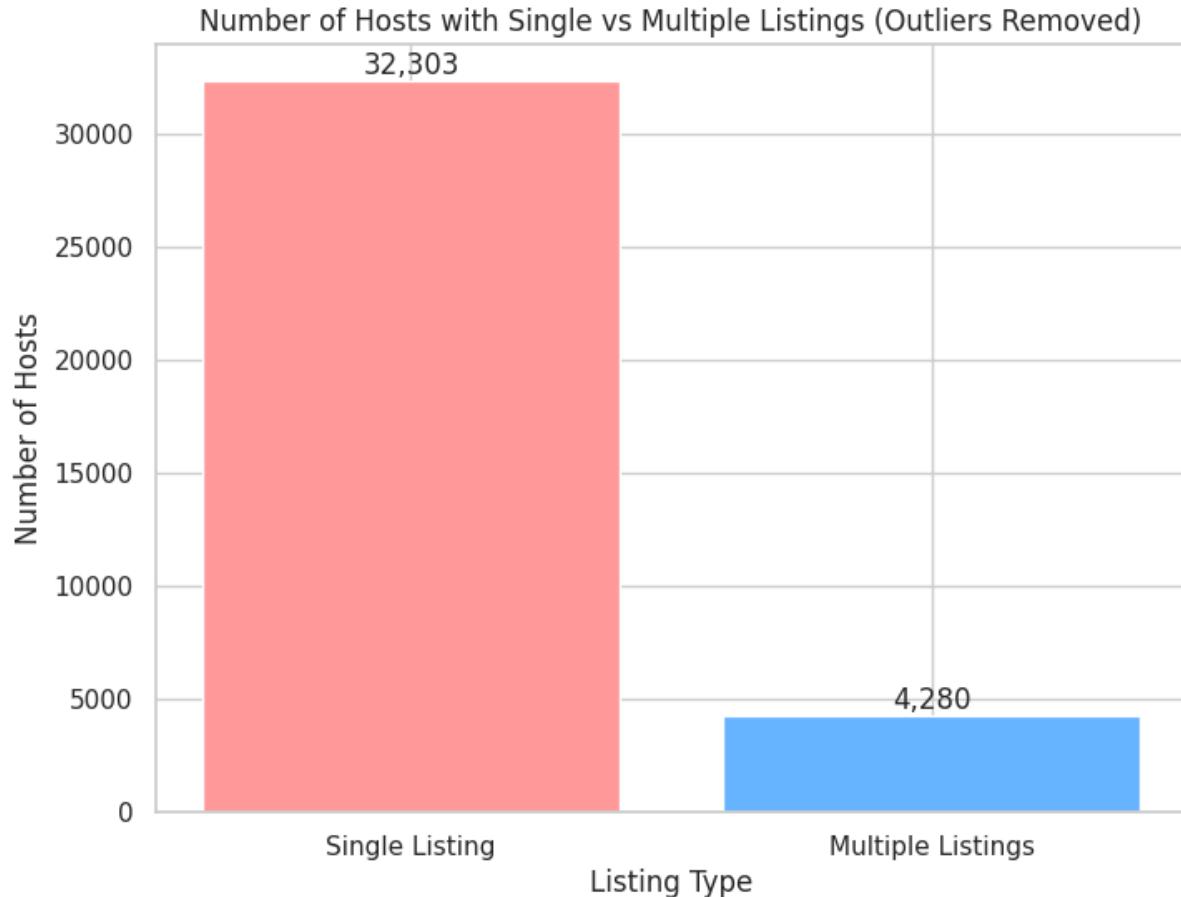
# Bar chart data
labels = ['Single Listing', 'Multiple Listings']
counts = [single_listing_count_filtered, multiple_listings_count_filtered]

# Create bar chart
plt.figure(figsize=(8, 6)) # Adjust figure size as needed
plt.bar(labels, counts, color=['#ff9999', '#66b3ff'])
plt.title('Number of Hosts with Single vs Multiple Listings (Outliers Removed)')
plt.xlabel('Listing Type')
plt.ylabel('Number of Hosts')

# Add count labels on top of bars
for i, v in enumerate(counts):
    plt.text(i, v + 2, f'{v:,}', ha='center', va='bottom') # Adjust positioning as needed

plt.show()

```



**Price (STEP 1)**  
**Statistical Descriptive Analysis on Price**  
In [35]:

```

# Mean price
print(f"Mean price: ${nyc_airbnb_df['price'].mean():.2f}")

# Median price
print(f"Median price: ${nyc_airbnb_df['price'].median():.2f}")

# Price range

```

```

print(f"Price range: ${nyc_airbnb_df['price'].min()} to ${nyc_airbnb_df['price'].max()}")
# Calculate quartiles
Q1 = nyc_airbnb_df['price'].quantile(0.25)
print(f"First Quartile (Q1) of price: ${Q1:.2f}")

Q3 = nyc_airbnb_df['price'].quantile(0.75)
print(f"Third Quartile (Q3) of price: ${Q3:.2f}")

# Calculate IQR
IQR = Q3 - Q1
print(f"Interquartile Range (IQR) of price: ${IQR:.2f}")

# Calculate outlier bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
print(f"Lower bound of price outliers: ${lower_bound:.2f}")
print(f"Upper bound of price outliers: ${upper_bound:.2f}")

# Number of outliers
outlier_below = nyc_airbnb_df[nyc_airbnb_df['price'] < (Q1 - 1.5 * IQR)]
outlier_above = nyc_airbnb_df[nyc_airbnb_df['price'] > (Q3 + 1.5 * IQR)]
print(f"Number of outliers below the IQR: {len(outlier_below)}")
print(f"Number of outliers above the IQR: {len(outlier_above)}")

# Total number of non-null price rows
non_null_price_count = nyc_airbnb_df['price'].notnull().sum()
print("Out of total number of non-null price rows:", non_null_price_count)

```

Mean price: \$152.72  
 Median price: \$106.00  
 Price range: \$0 to \$10000  
 First Quartile (Q1) of price: \$69.00  
 Third Quartile (Q3) of price: \$175.00  
 Interquartile Range (IQR) of price: \$106.00  
 Lower bound of price outliers: \$-90.00  
 Upper bound of price outliers: \$334.00  
 Number of outliers below the IQR: 0  
 Number of outliers above the IQR: 2972  
 Out of total number of non-null price rows: 48895  
 Review the high skew outliers (numeric).

- price
- minimum\_nights
- number\_of\_reviews
- reviews\_per\_month
- calculated\_host\_listings\_count

#### Price Box Plot

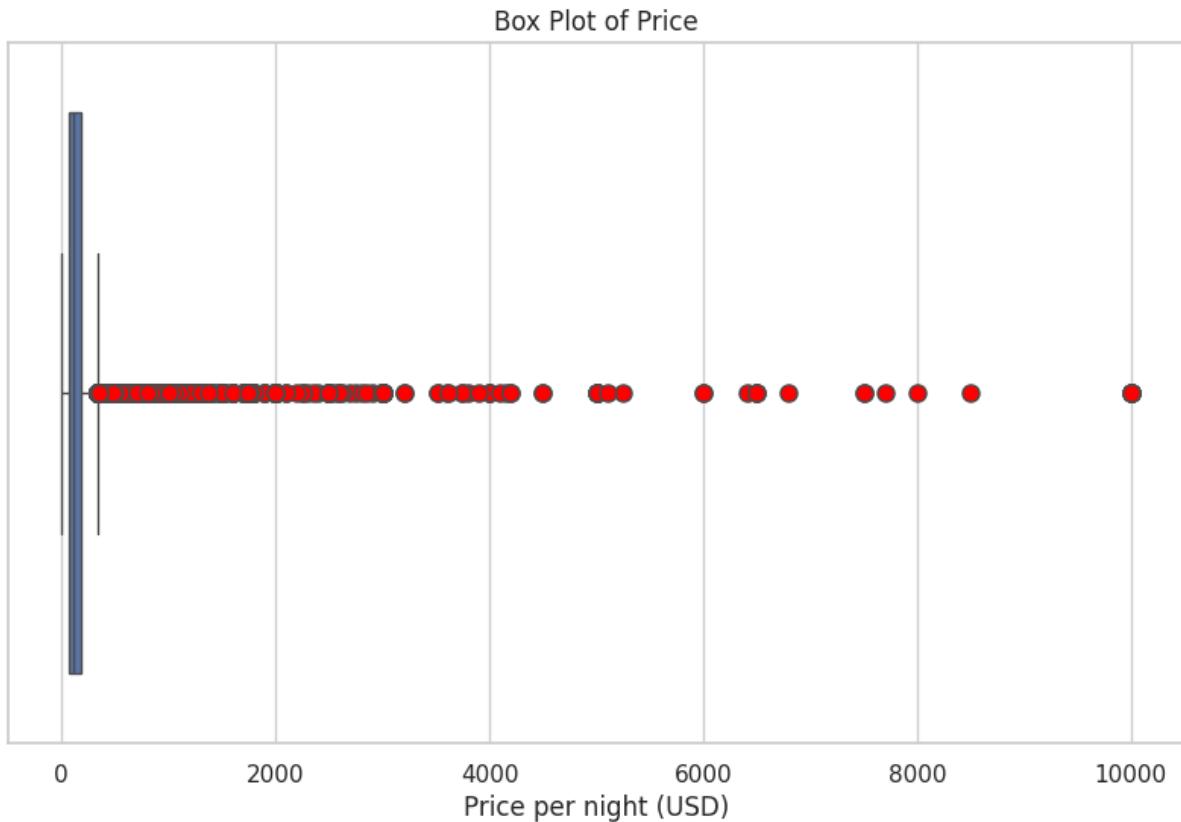
High priced listings could be luxury properties or in prime locations. Removing them will impact revenue. But they also skew average prices, making it harder to price new listings.  
 ACTION: Better understand what generates a higher price in bivariate. Is it a location, for instance?

In [36]:

```

# Create the box plot
plt.figure(figsize=(10, 6)) # Adjust figure size if needed
sns.boxplot(x='price', data=nyc_airbnb_df, # Use x='price' for horizontal orientation
            flierprops={'marker': 'o', 'markerfacecolor': 'red', 'markersize': 8}) # Customize outlier markers
plt.title('Box Plot of Price')
plt.xlabel('Price per night (USD)') # Change ylabel to xlabel
plt.show()

```



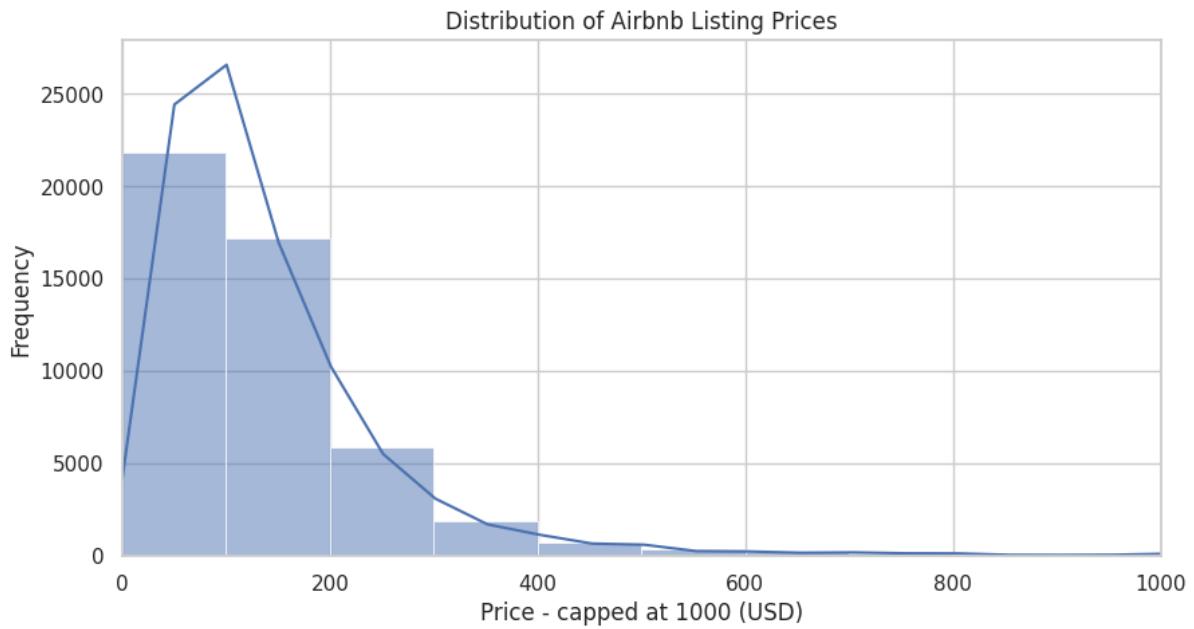
**Price Distribution (Zoomed in to 0-1000 USD)**

As seen in the calculations above, 334 USD is the upper bound above which there are 2,972 price outliers out of 48,895 non-null prices.

The zoomed in histogram shows that most listings are indeed priced below 400 USD, with a high concentration at lower prices.

In [37]:

```
# Plot the distribution of 'price' with a histogram capped at 1000 USD
plt.figure(figsize=(10, 5))
sns.histplot(nyc_airbnb_df['price'], bins=100, kde=True)
plt.xlim(0, 1000) # Limit x-axis for a clearer view (removes extreme outliers)
plt.title("Distribution of Airbnb Listing Prices")
plt.xlabel("Price - capped at 1000 (USD)")
plt.ylabel("Frequency")
plt.show()
```



#### Minimum nights box plot

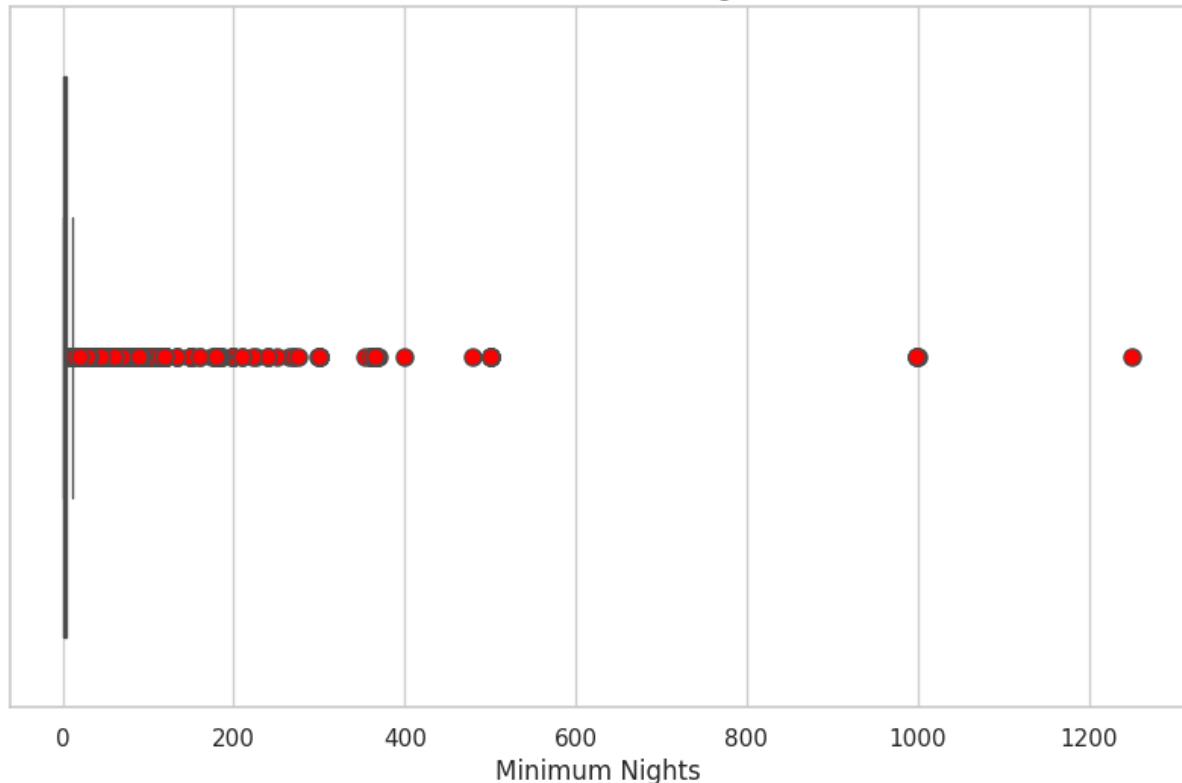
Very high minimum nights could be targeting specific customer segments like long-term rental.

ACTION: May want to split this into rental periods.

In [38]:

```
# Create the box plot
plt.figure(figsize=(10, 6)) # Adjust figure size if needed
sns.boxplot(x='minimum_nights', data=nyc_airbnb_df, # Use x='price' for horizontal orientation
            flierprops={'marker': 'o', 'markerfacecolor': 'red', 'markersize': 8}) # Customize outlier markers
plt.title('Box Plot of Minimum Nights')
plt.xlabel('Minimum Nights') # Change ylabel to xlabel
plt.show()
```

Box Plot of Minimum Nights



#### Number of Reviews

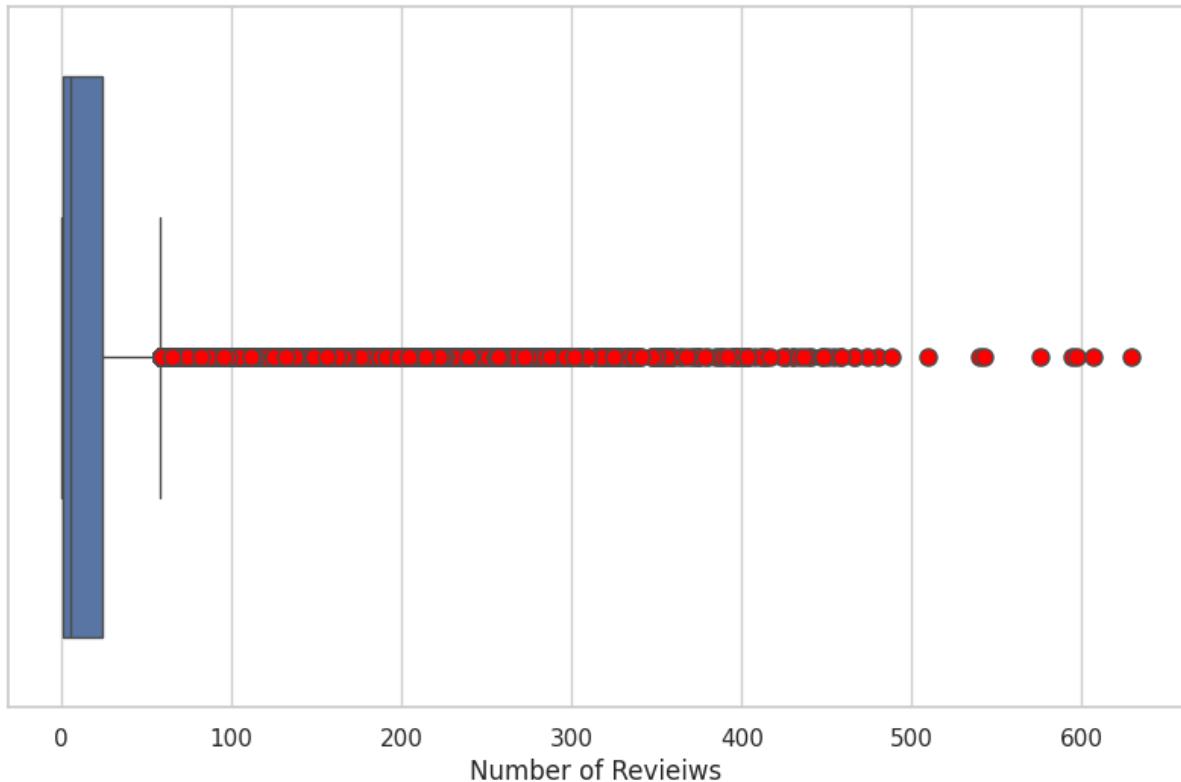
High number of reviews could mean popular listings or professional hosts (number listings) or rental for fewer days (minimum nights).

ACTION: Could compare to number of listings (i.e. professional hosts?), or see if more reviews have higher or lower prices, or see if it's just related to minimum nights (more frequently rented perhaps).

In [39]:

```
# Create the box plot
plt.figure(figsize=(10, 6)) # Adjust figure size if needed
sns.boxplot(x='number_of_reviews', data=nyc_airbnb_df, # Use x='price' for horizontal
            orientation
            flierprops={'marker': 'o', 'markerfacecolor': 'red', 'markersize': 8}) # Customize
            outlier markers
plt.title('Box Plot of Number of Reviews')
plt.xlabel('Number of Reviews') # Change ylabel to xlabel
plt.show()
```

Box Plot of Number of Reviews



#### Reviews per month

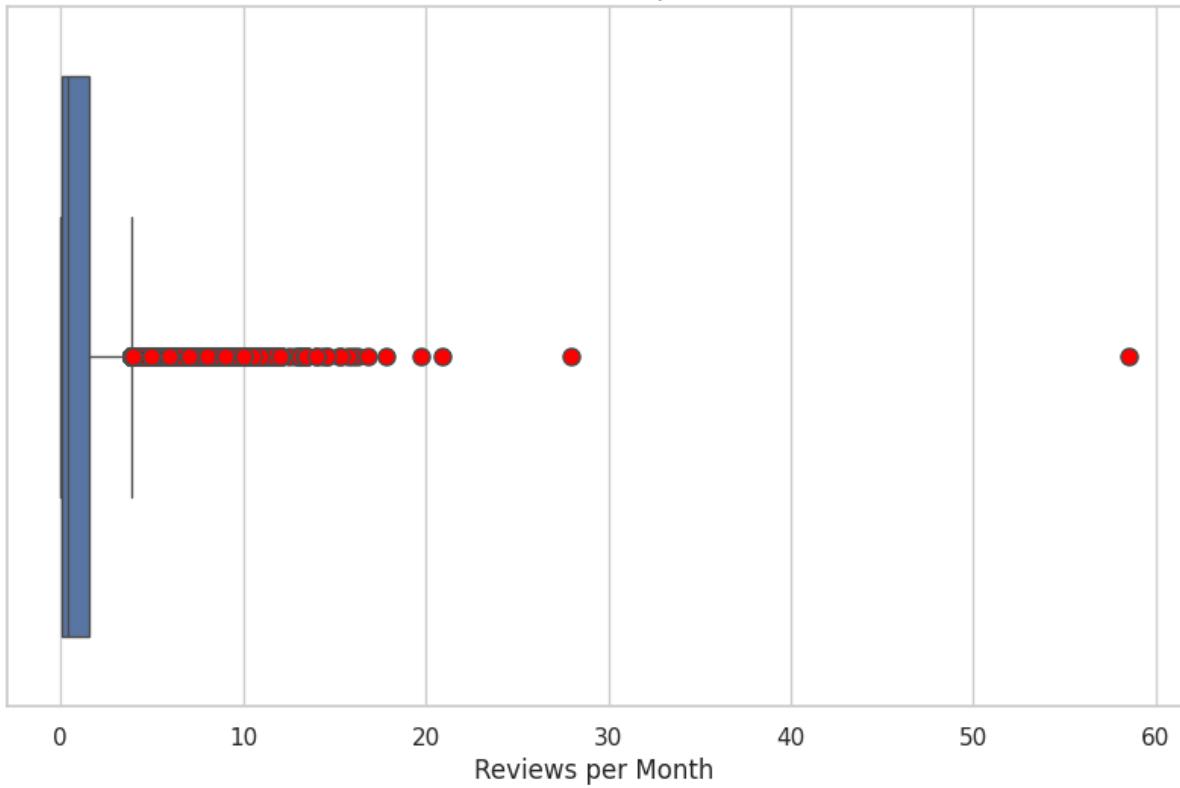
High number of reviews per month could mean popular listings or professional hosts (number listings) or rental for fewer days (minimum nights).

ACTION: Could compare to number of listings (i.e. professional hosts?), or see if more reviews have higher or lower prices, or see if it's just related to minimum nights (more frequently rented perhaps).

In [40]:

```
# Create the box plot
plt.figure(figsize=(10, 6)) # Adjust figure size if needed
sns.boxplot(x='reviews_per_month', data=nyc_airbnb_df, # Use x='price' for horizontal
            orientation
            flierprops={'marker': 'o', 'markerfacecolor': 'red', 'markersize': 8}) # Customize
            outlier markers
plt.title('Box Plot of Reviews per Month')
plt.xlabel('Reviews per Month') # Change ylabel to xlabel
plt.show()
```

Box Plot of Reviews per Month

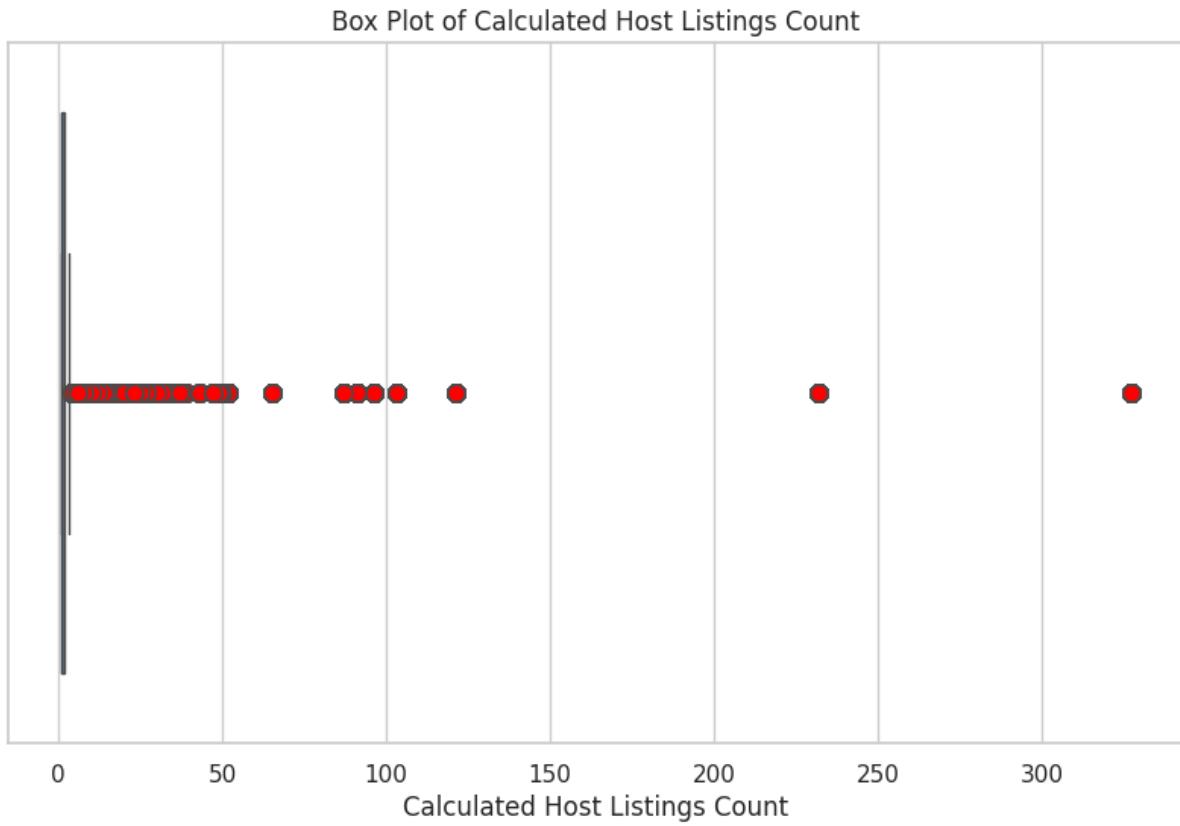


**Calculated Host Listings Count**

High number of listings probably means professional hosts with multiple properties.

ACTION: Could get a better idea of how many in each group to see how many professional hosts.  
In [41]:

```
# Create the box plot
plt.figure(figsize=(10, 6)) # Adjust figure size if needed
sns.boxplot(x='calculated_host_listings_count', data=nyc_airbnb_df, # Use
            x='calculated_host_listings_count' for horizontal orientation
            flierprops={'marker': 'o', 'markerfacecolor': 'red', 'markersize': 8}) # Customize
outlier markers
plt.title('Box Plot of Calculated Host Listings Count')
plt.xlabel('Calculated Host Listings Count') # Change ylabel to xlabel
plt.show()
```



#### BIVARIATE ANALYSIS

Oluveye (2023) for analysing two variables:

- crosstab / two-way tables
- pivot table
- scatter plot
- bar chart
- correlation analysis
- pairplots
- boxplots
- histograms

#### Correlation Heatmap

A correlation heatmap for the numeric columns helps understand their relationships better.

- +1: Perfect positive correlation
- +0.5 to +1: Strong positive correlation
- 0 to +0.5: Weak positive correlation
- 0: No correlation
- -0.5 to 0: Weak negative correlation
- -1 to -0.5: Weak negative correlation = -1: Perfect negative correlation

Potential for correlation:

- Could ID and number of reviews be how long someone has been a host?
- ID and host\_id are related positively, so there seems to be a relationship in ID creation.
- Number of reviews and reviews per month are positively related that makes sense.

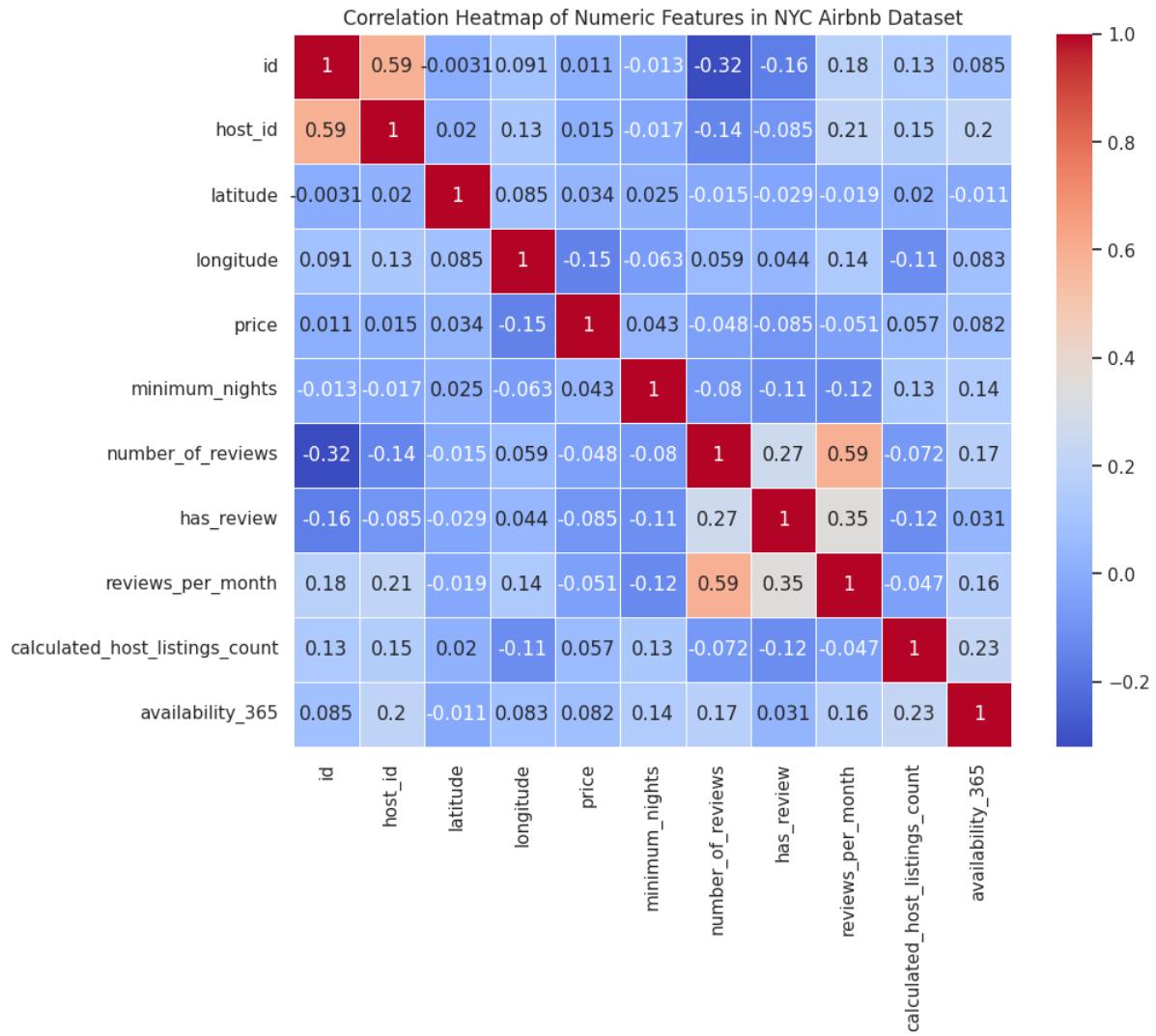
The correlation heatmap indicates the following relationships:

- Latitude and Longitude: They are not correlated with other variables in a meaningful way.
- Calculated Host Listings Count and Availability: Low correlation, suggesting that hosts's availability does not strongly correlate with their number of listings.
- Price and Minimum Nights: Surprisingly, there is no strong correlation between price and minimum\_nights, which may be further investigated for insights.

In [42]:

```
# Calculate correlation matrix
correlation_matrix_nyc = nyc_airbnb_df[numeric_columns_nyc].corr()
```

```
# Plot the heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix_nyc, annot=True, cmap="coolwarm", linewidths=0.5)
plt.title("Correlation Heatmap of Numeric Features in NYC Airbnb Dataset")
plt.show()
```



**Price by Location (STEP 2)**  
**Neighbourhood Group Versus Price (Zoomed Out)**  
In [43]:

```
nyc_airbnb_df.plot(kind='scatter', x='neighbourhood_group', y='price', s=32, alpha=.8)
plt.gca().spines[['top', 'right']].set_visible(False)
plt.title("Airbnb Listing Prices by Neighbourhood Group")
plt.xlabel("Neighbourhood Group")
plt.ylabel("Price (USD)")
plt.show()
```



#### **Airbnb Listing Prices by Neighbourhood Group - Zoomed In**

Examine how prices vary across different neighbourhood groups to see if location significantly influences price.

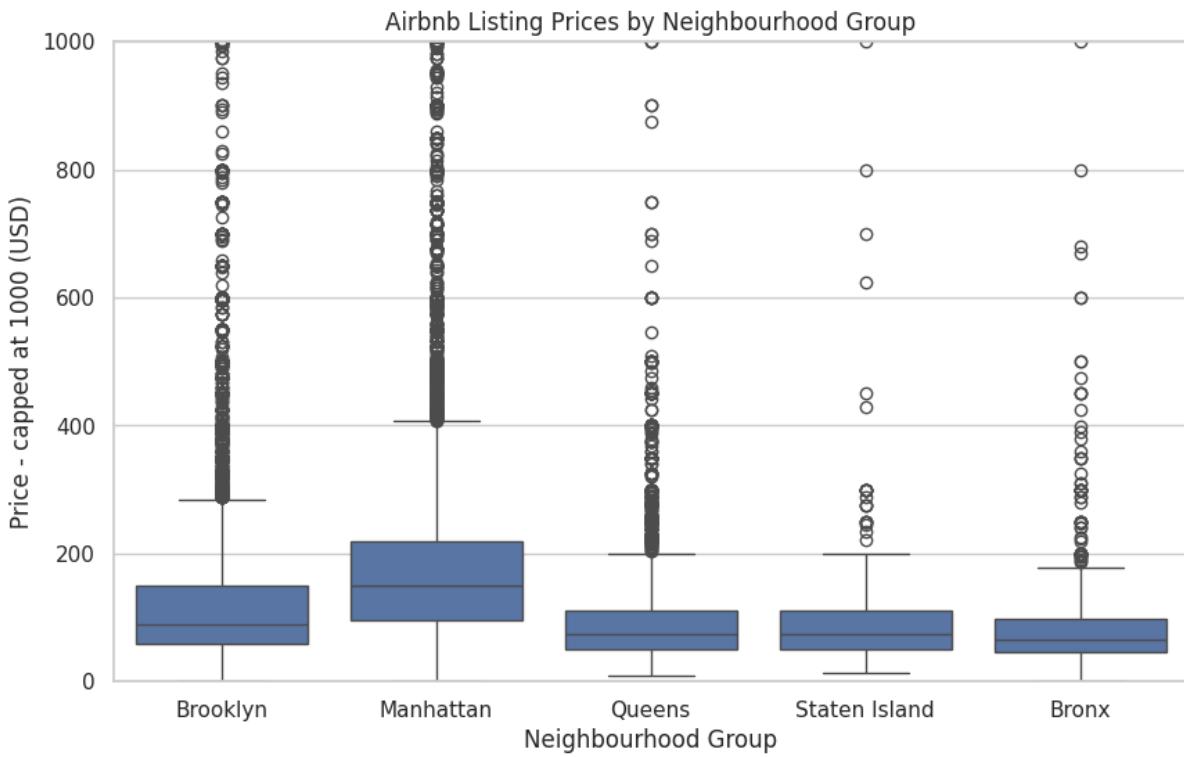
The boxplot indicates the following trends:

- Manhattan listings generally have the highest median prices, with more variability and outliers at the high end.
- Brooklyn follows Manhattan with moderately high prices. Queens, the Bronx, and Staten Island have lower median prices, with fewer outliers.
- This suggests that location (specifically, neighborhood) influences number of listings and prices.

#### **Price by Neighbourhood Group (Zoomed in 0-1000 USD)**

In [44]:

```
# Boxplot to show price distribution by neighbourhood group
plt.figure(figsize=(10, 6))
sns.boxplot(x='neighbourhood_group', y='price', data=nyc_airbnb_df)
plt.ylim(0, 1000) # Limit y-axis for a clearer view (removes extreme outliers)
plt.title("Airbnb Listing Prices by Neighbourhood Group")
plt.xlabel("Neighbourhood Group")
plt.ylabel("Price - capped at 1000 (USD)")
plt.show()
```



#### Median Price by Neighbourhood by Neighbourhood Group

Median rather than mean as less susceptible to outliers.  
In [45]:

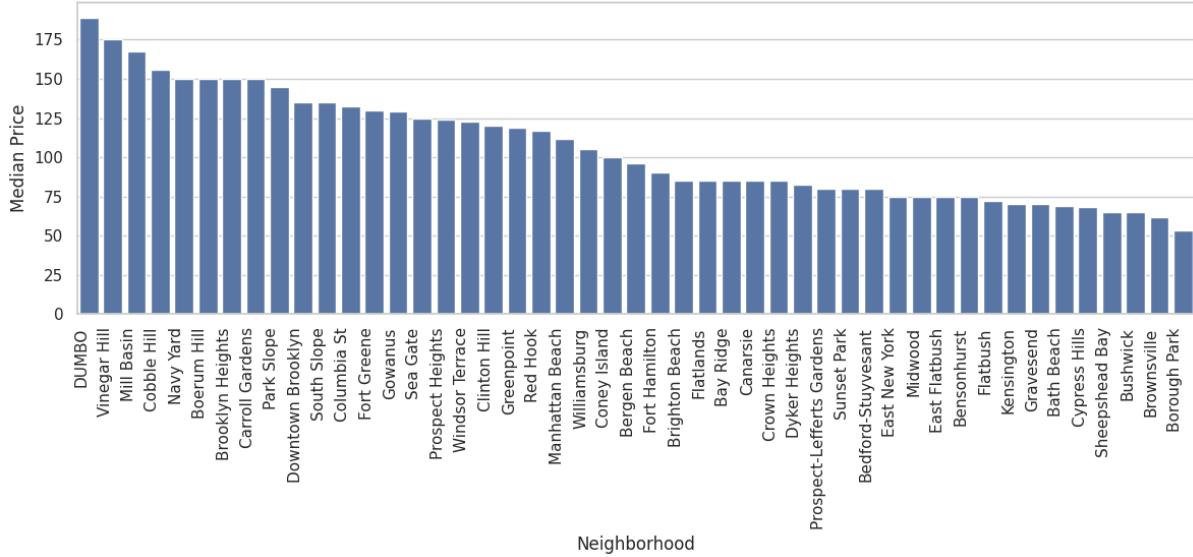
```
# Get unique neighborhood groups
neighborhood_groups = nyc_airbnb_df['neighbourhood_group'].unique()

# Loop through neighborhood groups and create plots
for group in neighborhood_groups:
    # Filter data for the current group
    group_data = nyc_airbnb_df[nyc_airbnb_df['neighbourhood_group'] == group]

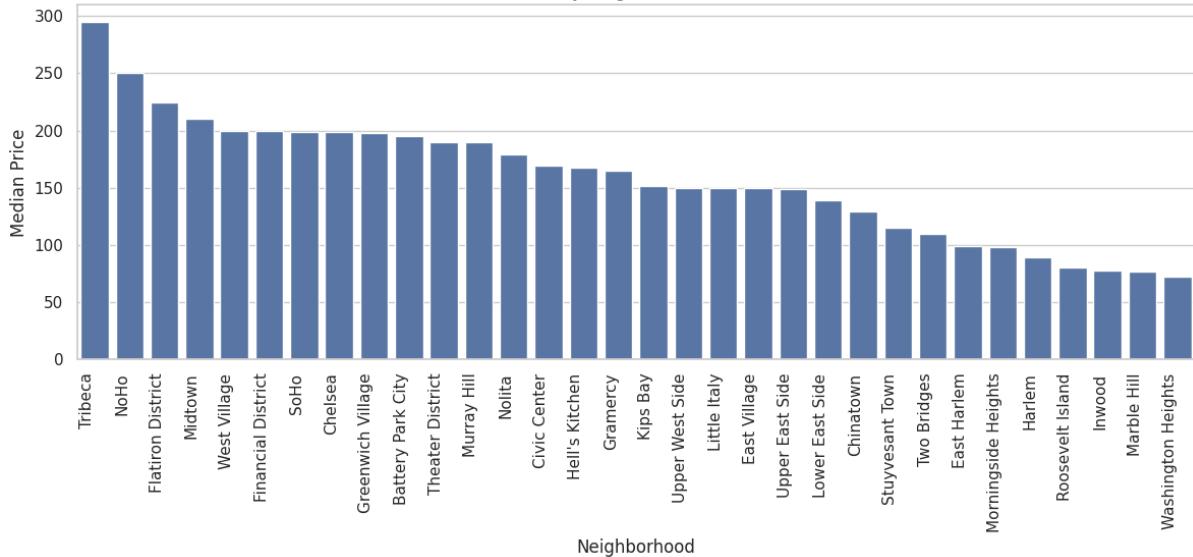
    # Calculate average price per neighborhood within the group
    median_price_by_neighborhood =
        group_data.groupby('neighbourhood')['price'].median().sort_values(ascending=False)

    # Create the bar chart
    plt.figure(figsize=(12, 6))  # Adjust figure size as needed
    sns.barplot(x=median_price_by_neighborhood.index, y=median_price_by_neighborhood.values)
    plt.title(f'Median Price by Neighborhood in {group}')
    plt.xlabel('Neighborhood')
    plt.ylabel('Median Price')
    plt.xticks(rotation=90, ha='right')  # Rotate x-axis labels for better readability
    plt.tight_layout()
    plt.show()
```

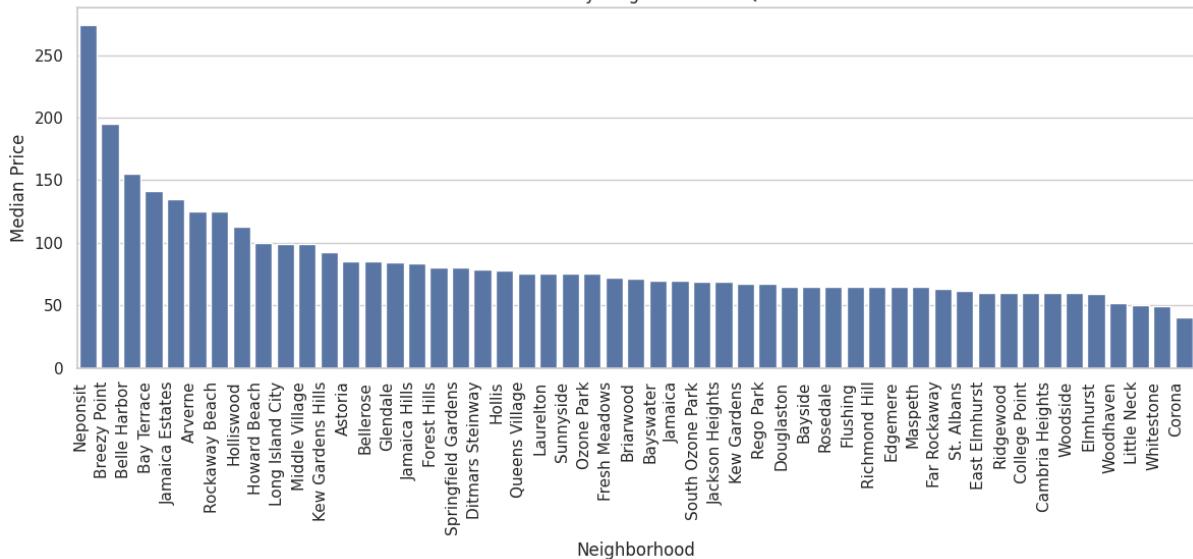
Median Price by Neighborhood in Brooklyn

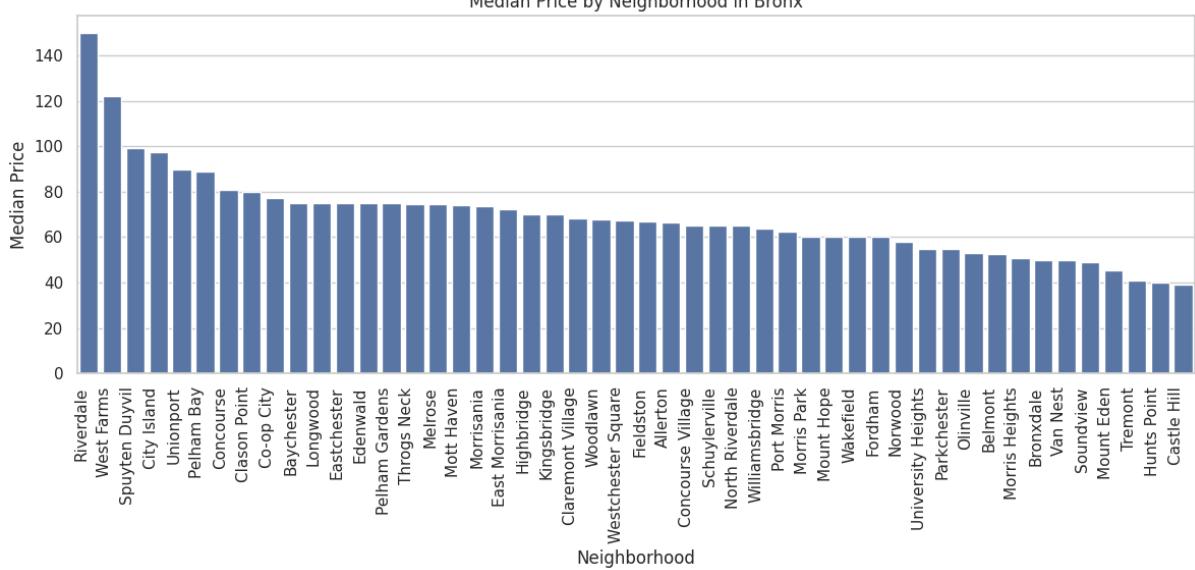
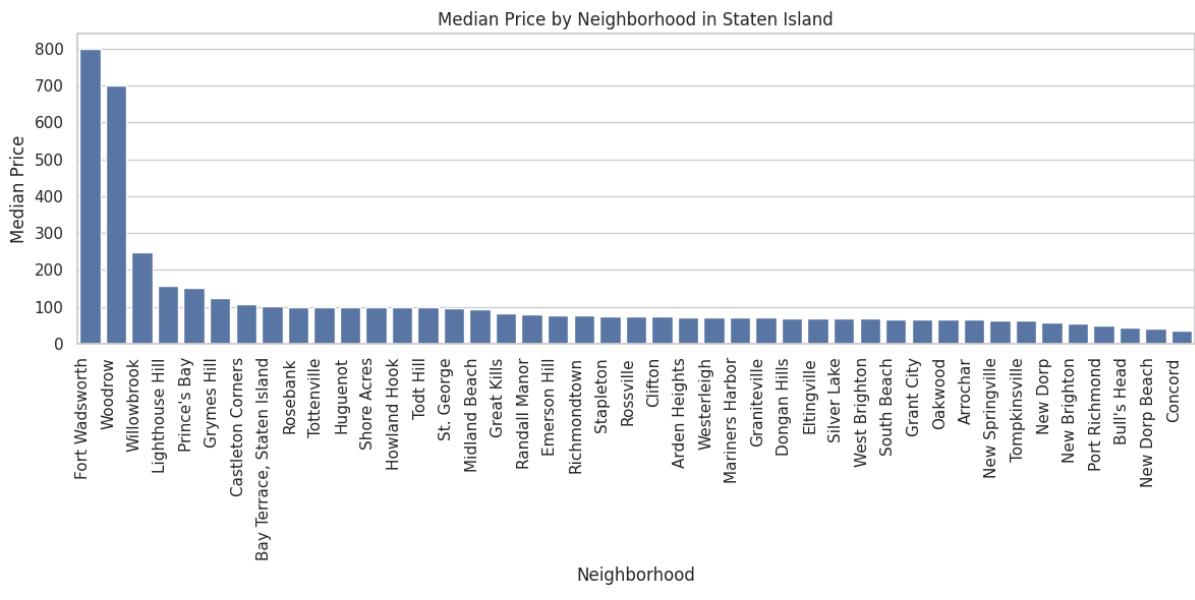


Median Price by Neighborhood in Manhattan



Median Price by Neighborhood in Queens





#### Top 3 highest priced neighbourhoods per neighbourhood group

What are the most desirable neighbourhoods? Are they real, outliers or just numbers thrown in as holding numbers? Could do further literary research on popular areas in NYC. (Not coding, but articles.)

In [46]:

```
def top_3_highest_price_neighborhoods(df):
    """
    Finds the top 3 neighborhoods with the highest individual listing prices
    within each neighborhood group, and the number of listings at that price.
    """

    Args:
        df: The input DataFrame containing Airbnb listing data.

    Returns:
        A DataFrame with the top 3 neighborhoods, their highest prices, and the
        number of listings at that price for each neighborhood group.
    """

    top_neighborhoods = []
    for group in df['neighbourhood_group'].unique():
        group_data = df[df['neighbourhood_group'] == group]

        # Sort by price in descending order and get the top 3 unique prices
        top_3_prices = group_data['price'].sort_values(ascending=False).unique()[:3]
```

```

for price in top_3_prices:
    # Find neighborhoods with listings at that price
    neighborhoods_at_price = group_data[group_data['price'] == price]['neighbourhood'].unique()

    for neighborhood in neighborhoods_at_price:
        # Count listings at that price in the neighborhood
        count = group_data[(group_data['price'] == price) &
(group_data['neighbourhood'] == neighborhood)].shape[0]
        top_neighborhoods.append([group, neighborhood, price, count])

return pd.DataFrame(top_neighborhoods, columns=['Neighborhood Group', 'Neighborhood',
'Highest Price', 'Count'])

# Get the top 3 neighborhoods by highest price with counts
top_neighborhoods_df = top_3_highest_price_neighborhoods(nyc_airbnb_df)

# Print the results
print(top_neighborhoods_df)

   Neighborhood Group      Neighborhood  Highest Price  Count
0       Brooklyn     Greenpoint      10000       1
1       Brooklyn    Clinton Hill      8000       1
2       Brooklyn    East Flatbush      7500       1
3    Manhattan  Upper West Side      10000       1
4    Manhattan    East Harlem      9999       1
5    Manhattan  Lower East Side      9999       2
6    Manhattan      Tribeca      8500       1
7      Queens      Astoria      10000       1
8      Queens      Bayside      2600       1
9      Queens   Forest Hills      2350       1
10   Staten Island  Randall Manor      5000       1
11   Staten Island   Prince's Bay      1250       1
12   Staten Island      St. George      1000       1
13      Bronx      Riverdale      2500       1
14      Bronx     City Island      1000       1
15      Bronx      Riverdale      800        1

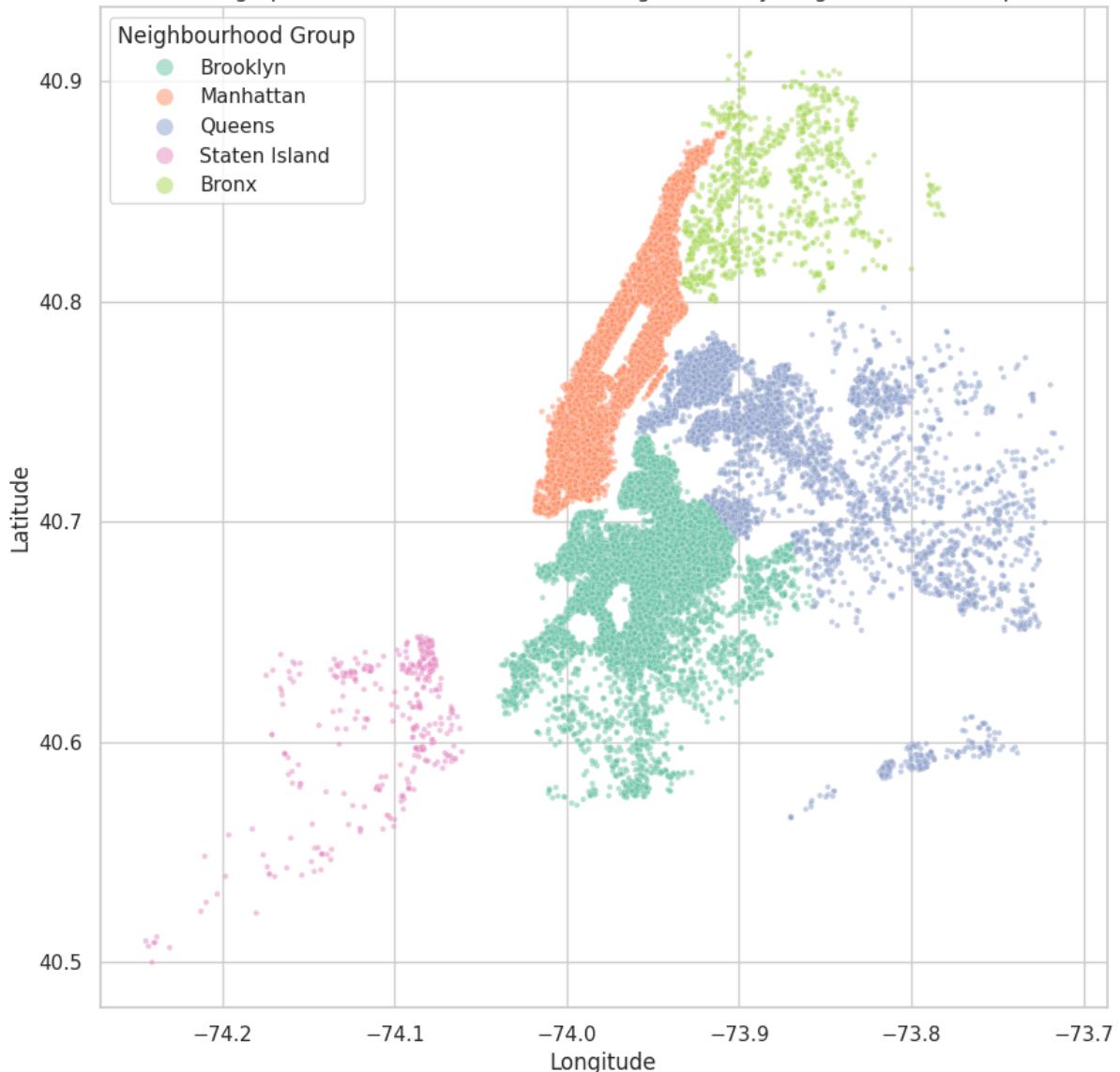
Geospatial Scatterplot
Basic geospatial analysis by plotting the locations of listings based on their latitude and longitude.
The geospatial plot reveals the distribution of Airbnb listings across New York City: Manhattan has a dense concentration of listings, especially in central and lower Manhattan. Brooklyn also has a significant number of listings, particularly in neighborhoods closer to Manhattan. Queens, the Bronx, and Staten Island have fewer listings, which are more spread out.
This map highlights the primary areas of Airbnb activity, aligning with the neighbourhoods in NYC.

In [47]:
```

```

# Geospatial analysis - plot listings by latitude and longitude
plt.figure(figsize=(10, 10))
sns.scatterplot(x='longitude', y='latitude', hue='neighbourhood_group', data=nyc_airbnb_df,
palette="Set2", s=10, alpha=0.5)
plt.title("Geographical Distribution of Airbnb Listings in NYC by Neighbourhood Group")
plt.xlabel("Longitude")
plt.ylabel("Latitude")
plt.legend(title='Neighbourhood Group', markerscale=3, scatterpoints=1)
plt.show()
```

Geographical Distribution of Airbnb Listings in NYC by Neighbourhood Group



In [48]:

```
plt.figure(figsize=(10, 10))
# Calculate normalized prices for alpha values beforehand
normalized_prices = nyc_airbnb_df['price'] / nyc_airbnb_df['price'].max()

# Get unique neighbourhood groups and their corresponding colors from the palette
neighbourhood_groups = nyc_airbnb_df['neighbourhood_group'].unique()
palette = sns.color_palette("Set2", n_colors=len(neighbourhood_groups))
group_colors = dict(zip(neighbourhood_groups, palette))

# Plot each neighbourhood group separately with varying transparency
for group in neighbourhood_groups:
    group_data = nyc_airbnb_df[nyc_airbnb_df['neighbourhood_group'] == group]
    plt.scatter(
        x=group_data['longitude'],
        y=group_data['latitude'],
        c=[group_colors[group]], # Use a single color for each group
        s=10,
        alpha=normalized_prices[group_data.index].values, # Use normalized prices for alpha
        label=group
    )

# Create custom legend handles with colors
legend_handles = [
    mlines.Line2D([], [], color=color, marker='o', linestyle='None', markersize=10,
    label=group)
```

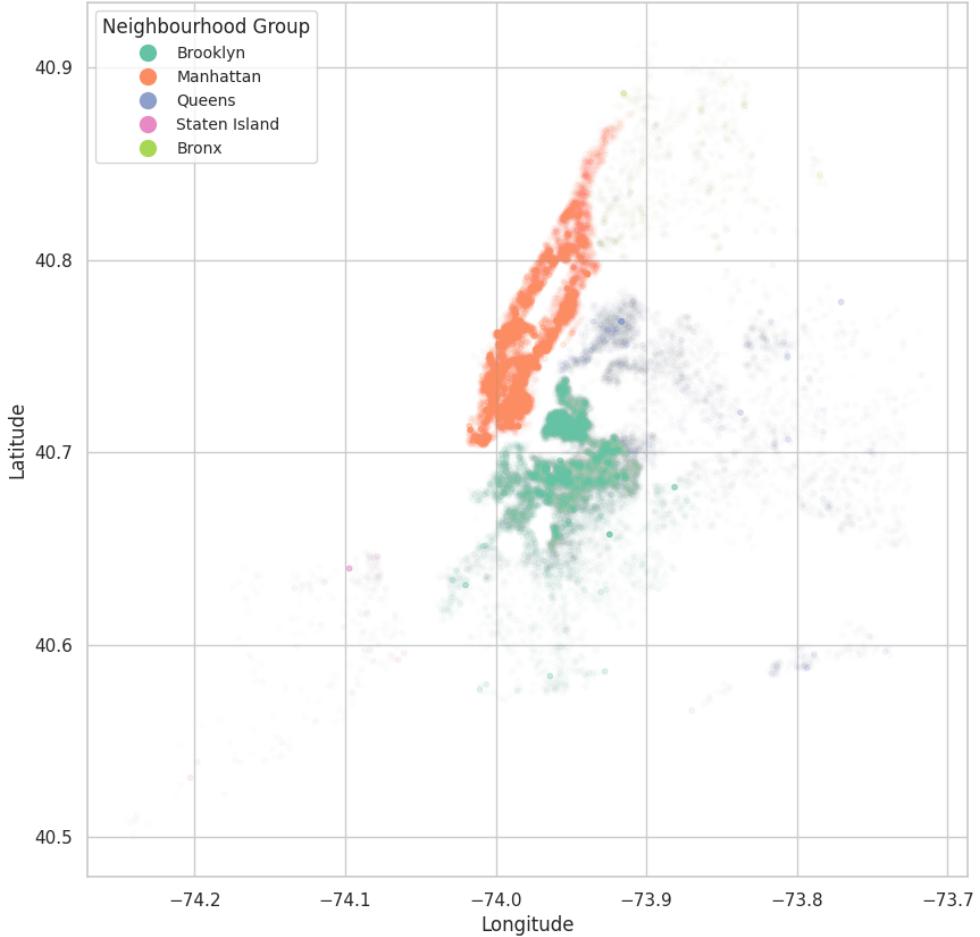
```

        for group, color in group_colors.items()
    ]

plt.title("Geographical Distribution of Airbnb Listings in NYC by Neighbourhood Group by Price Opacity", fontsize=16)
plt.xlabel("Longitude", fontsize=12)
plt.ylabel("Latitude", fontsize=12)
plt.legend(handles=legend_handles, title='Neighbourhood Group', fontsize=10) # Use custom handles
plt.show()

```

**Geographical Distribution of Airbnb Listings in NYC by Neighbourhood Group by Price Opacity**



#### Price by Room Type (STEP 3)

Next, analyze the price distribution by room type to assess if certain types of accommodations are priced higher than others

The boxplot reveals that:

- Entire home/apartment listings have the highest median prices, with a wide range and more high-priced outliers.
- Private rooms are generally priced lower, with less variability than entire homes.
- Shared rooms have the lowest median prices and minimal variability.

This suggests that room type impacts listing prices, with entire homes commanding the highest rates.

The scatter plot suggests that hosts with a larger number of listings generally set lower prices, with most high-priced listings concentrated among hosts with fewer listings. This could indicate that individual or small-scale hosts may price their properties higher, while larger hosts with multiple listings might offer competitive pricing.

#### Price by Room Type (Zoomed In)

Midline of box shows median (not mean) It is more robust to outliers, which makes it a better measure of central tendency than mean, which can be pulled higher.

In [49]:

```

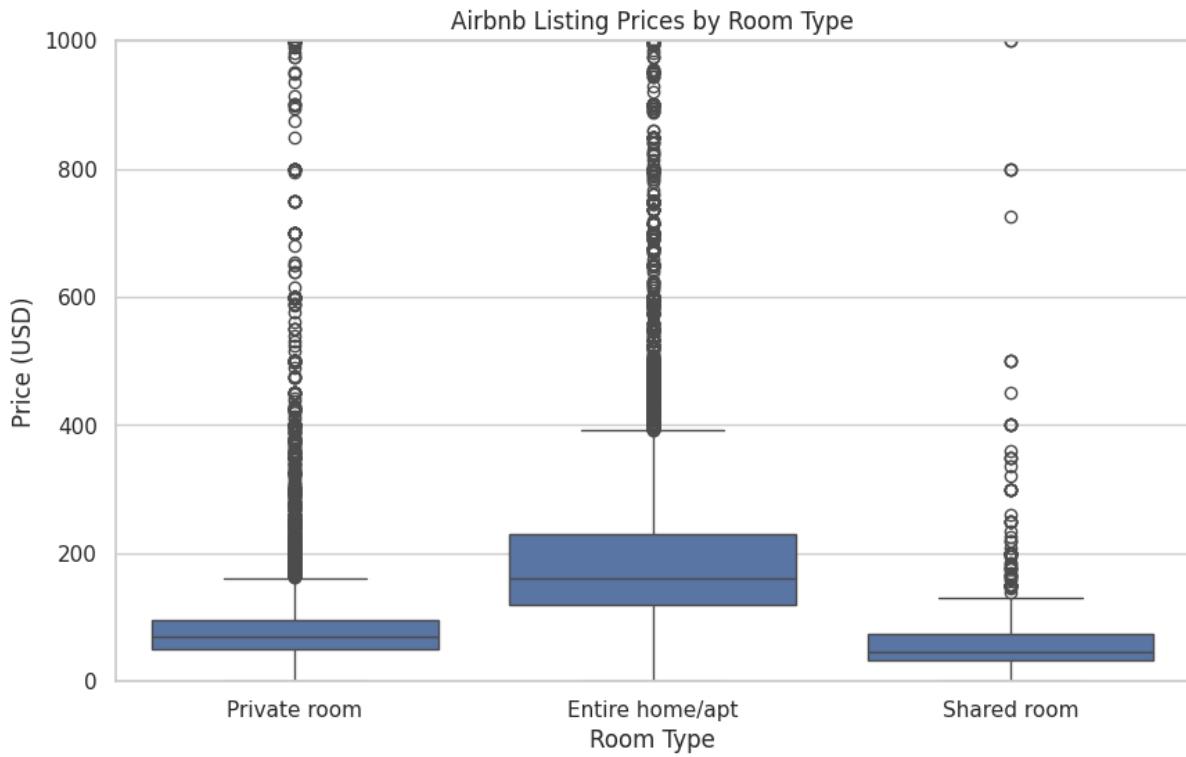
# Boxplot to show price distribution by room type
plt.figure(figsize=(10, 6))
sns.boxplot(x='room_type', y='price', data=nyc_airbnb_df)
plt.ylim(0, 1000) # Limit y-axis for a clearer view
plt.title("Airbnb Listing Prices by Room Type")

```

```

plt.xlabel("Room Type")
plt.ylabel("Price (USD)")
plt.show()

```



#### **Neighbourhood Group by Room Type by Price (STEP 5)**

Finally, examine the interaction between neighbourhood group and room type to see if combining these factors provides additional insights into pricing.

The boxplot highlights the combined influence of neighbourhood group and room type on pricing:

- Entire homes/apartments in Manhattan are priced significantly higher compared to other room types and neighborhoods, followed by Brooklyn.
- Private rooms show moderate pricing across neighborhoods, with Manhattan and Brooklyn remaining on the higher end.
- Shared rooms are generally the most affordable across all neighborhoods, with prices varying less between locations.

This interaction analysis reinforces that both neighborhood and room type play substantial roles in determining listing prices.

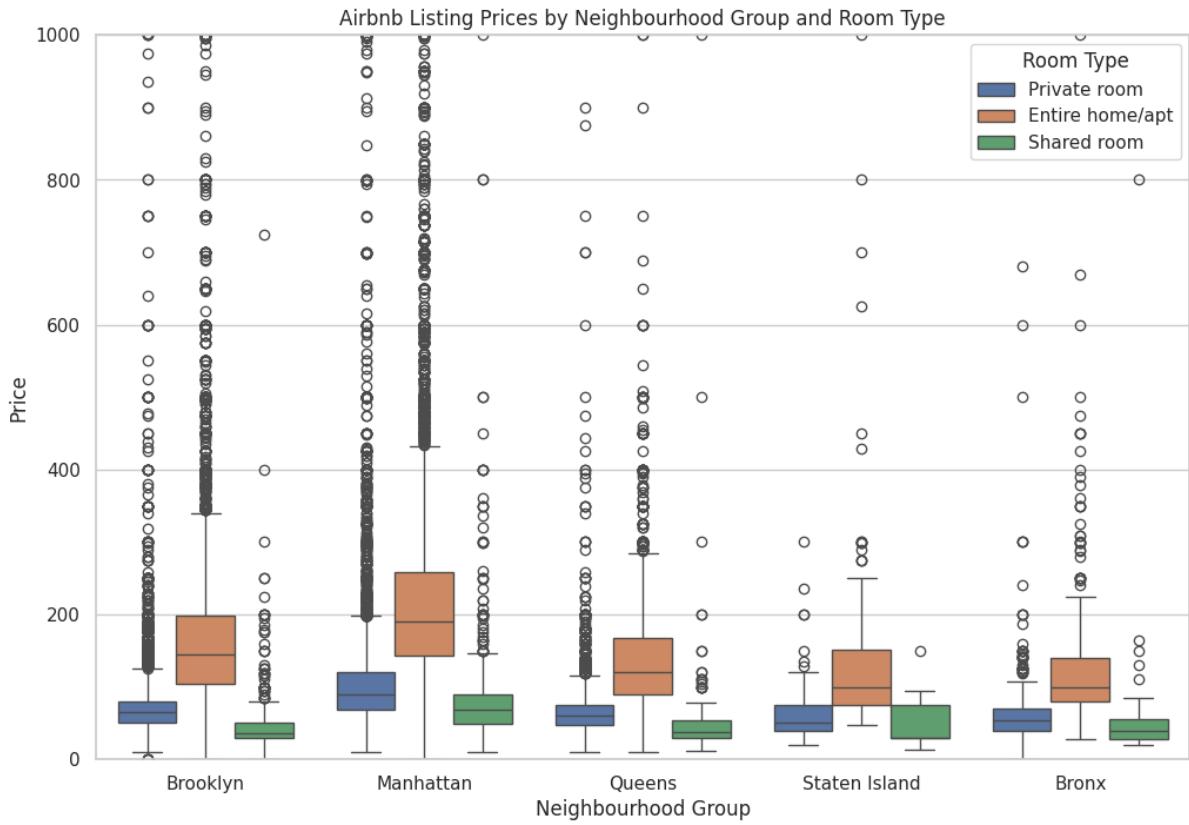
In [50]:

```
# Step 5: Price by Neighbourhood Group and Room Type Interaction
```

```

# Boxplot to show price distribution by neighbourhood group and room type
plt.figure(figsize=(12, 8))
sns.boxplot(x='neighbourhood_group', y='price', hue='room_type', data=nyc_airbnb_df)
plt.ylim(0, 1000) # Limit y-axis for clearer view
plt.title("Airbnb Listing Prices by Neighbourhood Group and Room Type")
plt.xlabel("Neighbourhood Group")
plt.ylabel("Price")
plt.legend(title="Room Type")
plt.show()

```



#### MULTIVARIATE ANALYSIS (MACHINE LEARNING)

Olueye (2023) describes multivariate using multiple variable techniques:

- k-means cluster analysis
- Principle Component Analysis (PCA)
- Factor analysis

Here, we only provide k-means cluster analysis.

#### Kmeans Cluster Analysis

##### Encode categorical features

In [51]:

```
# Encode and create mapping tables for categorical features
for col in ['room_type', 'neighbourhood_group']:
    encoder = LabelEncoder()
    nyc_airbnb_df[col + '_encoded'] = encoder.fit_transform(nyc_airbnb_df[col])
    mapping = dict(zip(encoder.classes_, encoder.transform(encoder.classes_)))
    print(f"{col.title()} Mapping:", mapping)
```

Room Type Mapping: {'Entire home/apt': 0, 'Private room': 1, 'Shared room': 2}  
 Neighbourhood\_Group Mapping: {'Bronx': 0, 'Brooklyn': 1, 'Manhattan': 2, 'Queens': 3, 'Staten Island': 4}

#### Feature Selection

In [52]:

```
# Define selected features for analysis
# Choose from: 'latitude', 'longitude', 'price', 'minimum_nights', 'availability_365',
# 'number_of_reviews', 'reviews_per_month', 'room_type_encoded', 'neighbourhood_group_encoded'
# features = ['price', 'minimum_nights', 'number_of_reviews', 'availability_365']
features = ['price', 'minimum_nights', 'availability_365', 'number_of_reviews',
            'reviews_per_month', 'calculated_host_listings_count', 'room_type_encoded',
            'neighbourhood_group_encoded']

# Drop rows with NaN values in the selected features
data_cleaned = nyc_airbnb_df.dropna(subset=features).copy()
```

#### Normality Test

- Can not use Shapiro-Wilk here because it is designed for data sets less than 5000.
- Use Anderson-Darling because large dataset with outliers.

Anderson-Darling shows:

- Choose 5.0% significance level 0.787

- Null hypothesis: If Anderson-Darling statistic is > critical value at 5%, reject null hypothesis of normality.
- Alternative: If Anderson-Darling is <= critical value at 5%, fail to reject null.
- all features are > 0.787 ==> Reject null, not normal.

As seen throughout, they are not normal.

In [53]:

```
# Normality test on the selected features
normality_results = {}
for col in features:
    # Drop NaN values before the test
    data_col = data_cleaned[col].dropna()

    # Perform Anderson-Darling test
    statistic, critical_values, significance_level = stats.anderson(data_col, dist='norm')

    # Store results in the dictionary
    normality_results[col] = {
        'Anderson-Darling Statistic': statistic,
        'Critical Values': critical_values,
        'Significance Level': significance_level
    }

# Convert the results to a DataFrame for better readability
normality_results_df = pd.DataFrame(normality_results).T
normality_results_df.reset_index(inplace=True)
normality_results_df.rename(columns={'index': 'Feature'}, inplace=True)

# Display the normality results in a clean table format
print("Normality Test Results (Anderson-Darling):")
print(normality_results_df)

Normality Test Results (Anderson-Darling):
                                         Feature Anderson-Darling Statistic          Critical
Values \
0                               price           7277.445657  [0.576, 0.656, 0.787, 0.918,
1.092]
1                         minimum_nights       10324.332785  [0.576, 0.656, 0.787, 0.918,
1.092]
2                     availability_365        4044.550982  [0.576, 0.656, 0.787, 0.918,
1.092]
3                  number_of_reviews       7077.730483  [0.576, 0.656, 0.787, 0.918,
1.092]
4                reviews_per_month        4416.952658  [0.576, 0.656, 0.787, 0.918,
1.092]
5  calculated_host_listings_count      15451.861681  [0.576, 0.656, 0.787, 0.918,
1.092]
6                 room_type_encoded       7335.071753  [0.576, 0.656, 0.787, 0.918,
1.092]
7            neighbourhood_group_encoded   3937.089193  [0.576, 0.656, 0.787, 0.918,
1.092]

                                         Significance Level
0  [15.0, 10.0, 5.0, 2.5, 1.0]
1  [15.0, 10.0, 5.0, 2.5, 1.0]
2  [15.0, 10.0, 5.0, 2.5, 1.0]
3  [15.0, 10.0, 5.0, 2.5, 1.0]
4  [15.0, 10.0, 5.0, 2.5, 1.0]
5  [15.0, 10.0, 5.0, 2.5, 1.0]
6  [15.0, 10.0, 5.0, 2.5, 1.0]
7  [15.0, 10.0, 5.0, 2.5, 1.0]

Remove Outliers
• k-means is sensitive to outliers.
• Remove outliers using Interquartile Range (IQR) method.
```

In [54]:

```
# Remove outliers using IQR and standardize the data
def remove_outliers_iqr(dataframe, columns):
    cleaned_data = dataframe.copy()
    for col in columns:
        Q1 = cleaned_data[col].quantile(0.25)
        Q3 = cleaned_data[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
```

```

        cleaned_data = cleaned_data[(cleaned_data[col] >= lower_bound) & (cleaned_data[col] <=
upper_bound)]
    return cleaned_data

# Applying the outlier removal
data_no_outliers = remove_outliers_iqr(data_cleaned, features)
data_no_outliers.info()

<class 'pandas.core.frame.DataFrame'>
Index: 29408 entries, 1 to 48894
Data columns (total 17 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   id               29408 non-null   int64  
 1   host_id          29408 non-null   int64  
 2   neighbourhood_group  29408 non-null   object  
 3   neighbourhood      29408 non-null   object  
 4   latitude          29408 non-null   float64 
 5   longitude         29408 non-null   float64 
 6   room_type         29408 non-null   object  
 7   price             29408 non-null   int64  
 8   minimum_nights    29408 non-null   int64  
 9   number_of_reviews  29408 non-null   int64  
 10  last_review       29408 non-null   object  
 11  has_review        29408 non-null   int64  
 12  reviews_per_month 29408 non-null   float64 
 13  calculated_host_listings_count 29408 non-null   int64  
 14  availability_365   29408 non-null   int64  
 15  room_type_encoded  29408 non-null   int64  
 16  neighbourhood_group_encoded 29408 non-null   int64  
dtypes: float64(3), int64(10), object(4)
memory usage: 4.0+ MB
In [55]:
```

# Display the data without outliers

```

data_no_outliers.reset_index(drop=True, inplace=True)  # Resetting index for clarity
data_no_outliers.head()  # Displaying the first few rows of the dataset without outliers
```

Out[55]:

			l	l	r	mi	num	ber	las	ha	rev	calcul	ava	roo	neigh
h	nei	ne	a	o	o	p	ni	ber	las	ha	iew	ated_h	ila	m_t	bourn
o	ghb	ig	t	n	o	r	mu	_of	t_r	s	s_p	ost_li	bil	ype	ood_g
i	our	hb	i	g	m	r	mu	_re	evi	re	er	ity	_en	roup_	
s	hoo	ou	t	i	i	i	m_	vie	ew	ew	mon	stings	_36	cod	encod
t	d_g	rh	t	t	t	c	ni	re	ew	ew	th	_count	5	ed	ed
—	rou	oo	u	u	y	e	gh	vie	ew	ew					
d	p	d	d	d	p	ts	ws								

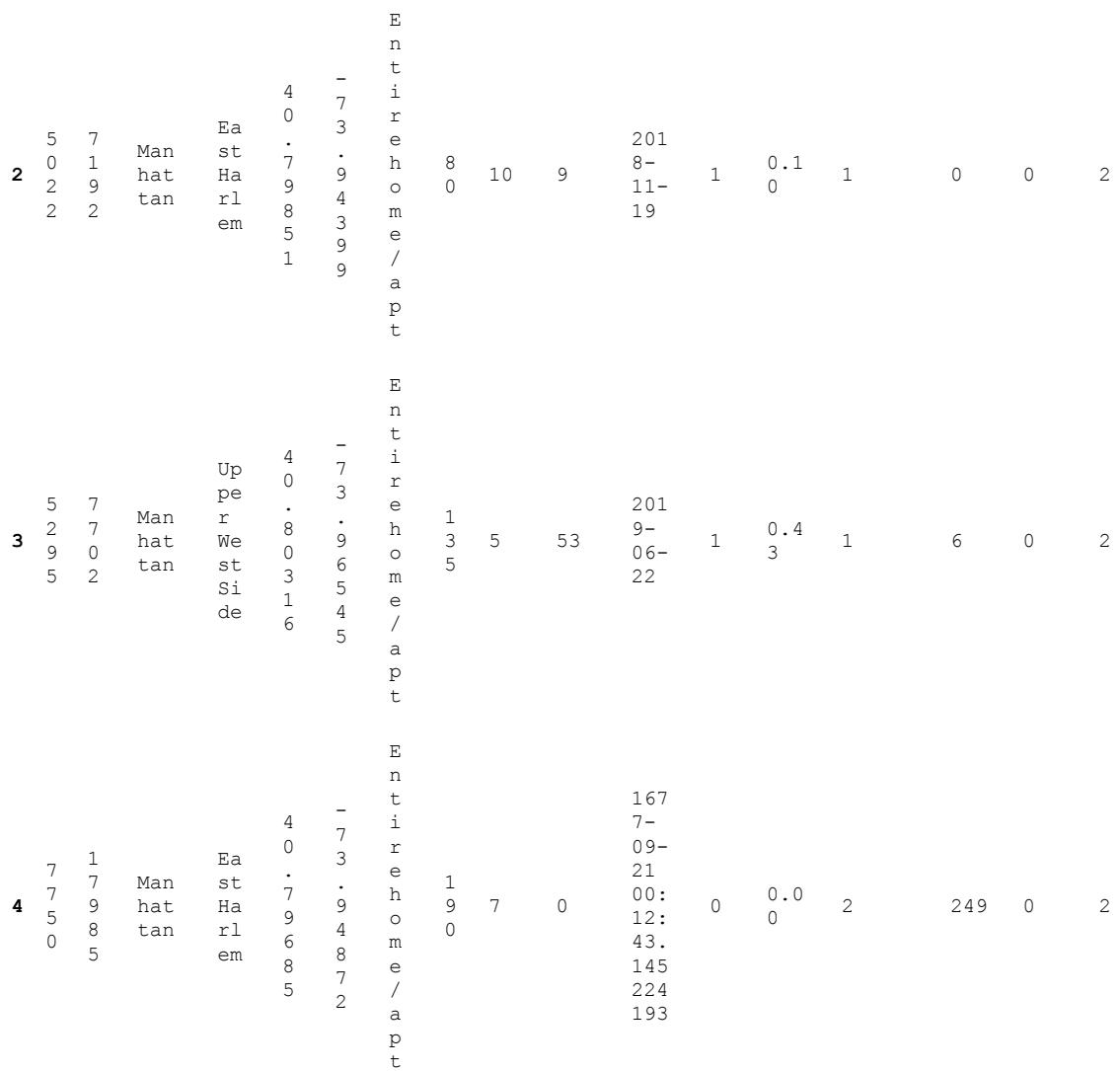
  

			E													
0	2	2	Man	Mi	.	e										
5	5	8	hat	dt	7	.	h	2								
9	9	4	hat	ow	5	9	o	2	1	45	9-					
5	5	5	tan	n	3	8	m	5			05-	1	0.3	2	355	0
											8					2

			P													
1	3	4	Man	Ha	.	-	r									167
6	6	6	hat	rl	8	.	a	1								7-
4	4	3	tan	em	0	9	t	5	3	0	00:	0	0.0	1	365	1
7	7	2			9	4	e	0			12:	0	0			2
					0	1	r				43.					
					2	9	o				145					
					0	0	o				224					
					2		m				193					

	h	nei	ne	l	o	r	mi	num	ha	rev	calcul	ava	roo	neigh	
o	ghb	ig	a	n	o	p	ni	ber	las	ha	iew	ilated_h	ila	m_t	bourn
s	our	hb	t	g	m	r	mu	_of	t_r	s_	s_p	ost_li	bil	ype	ood_g
i	t	hoo	ou	i	i	i	m_	_re	evi	re	er	stings	ity	_en	roup
d	d_g	rh	u	t	t	c	ni	vie	ew	vi	mon	_count	36	cod	encod
i	rou	oo	d	u	y	e	gh	ws	ew	th	th	5	ed	ed	
d	p	d	d	d	p	ts									



#### Linear Correlation Map without outliers

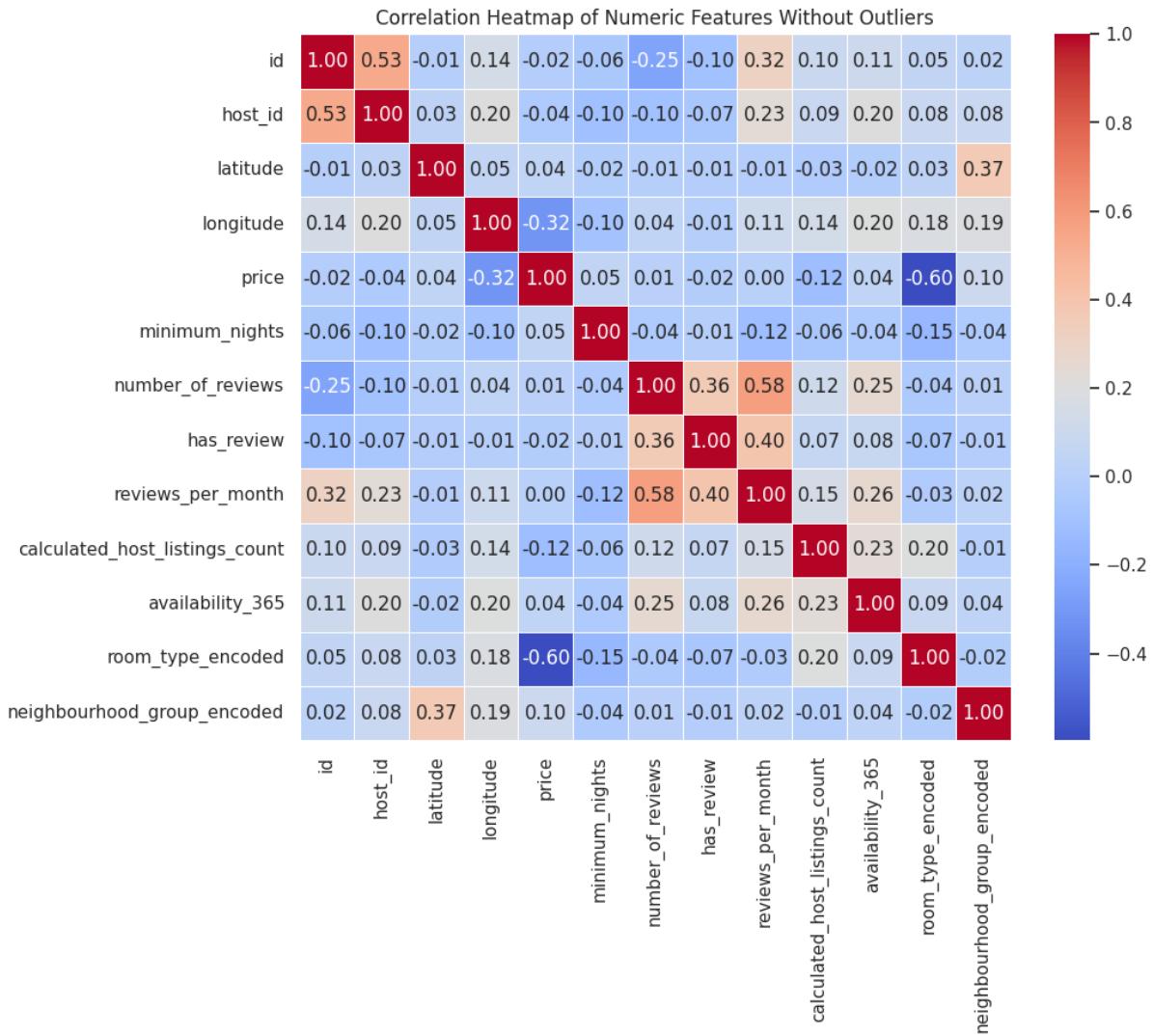
In [57]:

```
# Get new numeric list with encoded categoricals
numeric_columns_nyc = data_no_outliers.select_dtypes(include=['float64', 'int64']).columns

# Outlier free correlation matrix
correlation_matrix_no_outliers = data_no_outliers[numeric_columns_nyc].corr()

# Plot the heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix_no_outliers, annot=True, cmap="coolwarm", linewidths=0.5,
fmt=".2f")

plt.title("Correlation Heatmap of Numeric Features Without Outliers")
plt.show()
```



#### Log Correlation Map without outliers

In [58]:

```
# Create a copy of the DataFrame to avoid modifying the original
data_log_transformed = data_no_outliers.copy()

# Apply log transformations to selected columns
for col in ['price', 'minimum_nights', 'number_of_reviews', 'reviews_per_month',
'calculated_host_listings_count']:
    data_log_transformed[col + '_log'] = np.log1p(data_log_transformed[col])

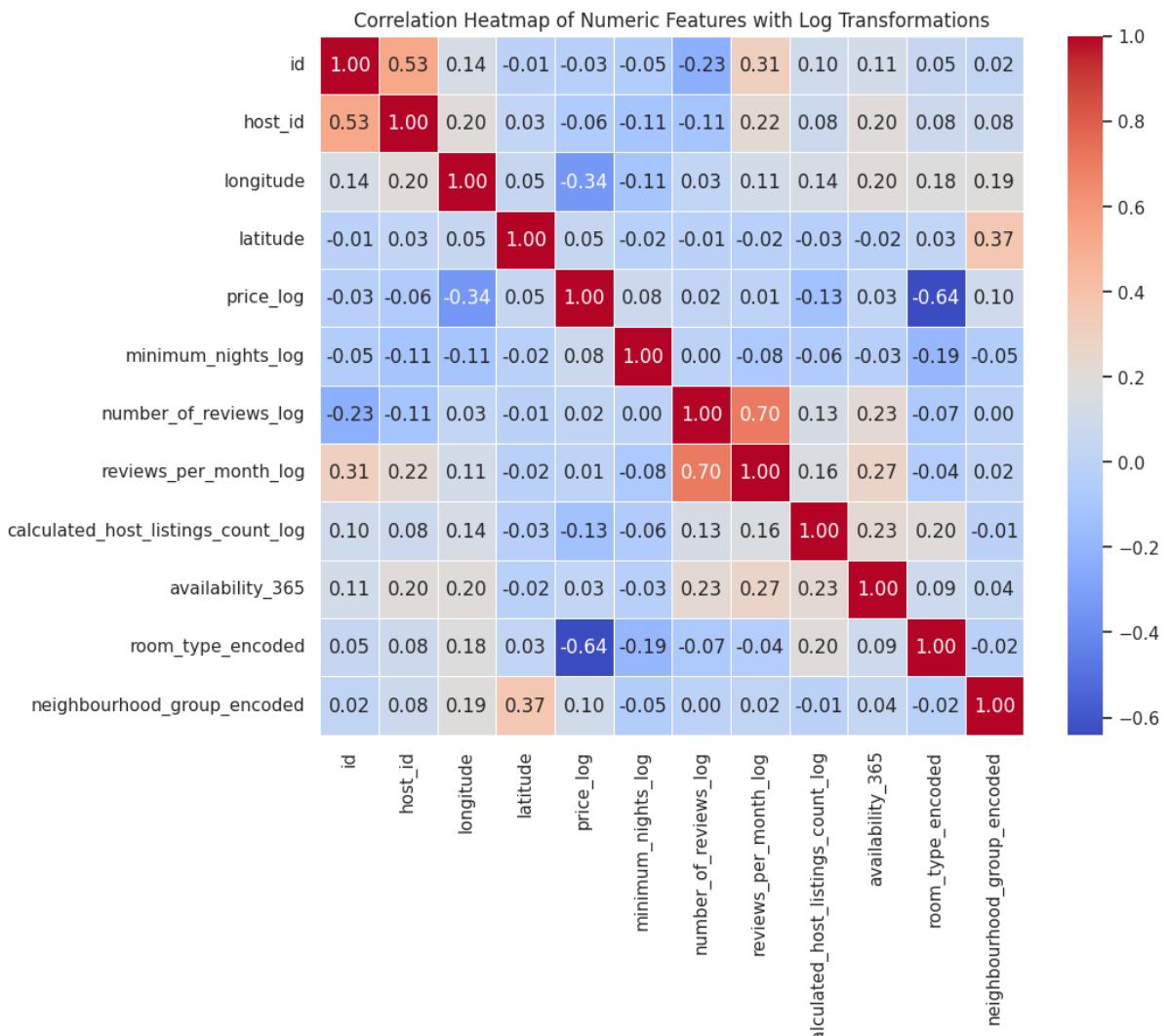
# Create a list of log-transformed features
log_transformed_features = ['price_log', 'minimum_nights_log', 'number_of_reviews_log',
                           'reviews_per_month_log', 'calculated_host_listings_count_log']

# Select columns to include and remove NaN values
data_for_correlation = data_log_transformed[['id', 'host_id', 'longitude', 'latitude'] +
log_transformed_features +
['availability_365', 'room_type_encoded', 'neighbourhood_group_encoded']].dropna()

# ... Update numeric_columns_nyc to include 'price_log' and others
numeric_columns_nyc = data_for_correlation.select_dtypes(include=['float64', 'int64']).columns

# Calculate correlation matrix
correlation_matrix_log_transformed = data_for_correlation[numeric_columns_nyc].corr()

# Plot the heatmap
plt.figure(figsize=(10, 8)) # Increased figure size to accommodate labels
sns.heatmap(correlation_matrix_log_transformed.round(2), annot=True, cmap="coolwarm",
linewdiths=0.5, fmt=".2f")
plt.title("Correlation Heatmap of Numeric Features with Log Transformations")
plt.show()
```



#### Standardise Data

In [59]:

```
# Standardizing the data
scaler = StandardScaler()
scaled_data_no_outliers = scaler.fit_transform(data_no_outliers[features])

print("Outlier removal and standardization are complete.")
```

Outlier removal and standardization are complete.

#### Display standardised data without outliers

In [60]:

```
# Display the standardised data without outliers
scaled_data_no_outliers_df = pd.DataFrame(scaled_data_no_outliers, columns=features)
scaled_data_no_outliers_df.head()
```

Out[60]:

	<code>pri ce</code>	<code>minimu m_nigh ts</code>	<code>availab ility_3 65</code>	<code>number_o f_review s</code>	<code>reviews_ per_mont h</code>	<code>calculated_hos t_listings_cou nt</code>	<code>room_typ e_encode d</code>	<code>neighbourhood group_encode d</code>
0	1.5 424 96	- 0.9726 69	2.50641 2	2.260172	- 0.271607	1.329224	- 0.975070	0.540583
1	0.4 270 50	0.0451 49	2.59463 4	- 0.691617	- 0.773451	-0.470823	0.887417	0.540583

	pri ce	minimu m_nigh ts	availab ility_3	number_o f_review s	reviews_ per_mont h	calculated Hos t_listings_cou nt	room_typ e_encode d	neighbourhood group_encode d
2	0.6 140 34	3.6075 11	65	- 0.62548 8	- 0.101259	- 0.641387	-0.470823	- 0.975070 0.540583
3	0.2 039 60	1.0629 67	- 0.57255 5	2.784934	- 0.205575	-0.470823	- 0.975070 0.540583	
4	1.0 219 54	2.0807 85	1.57125 3	- 0.691617	- 0.773451	1.329224	- 0.975070 0.540583	

#### Select k

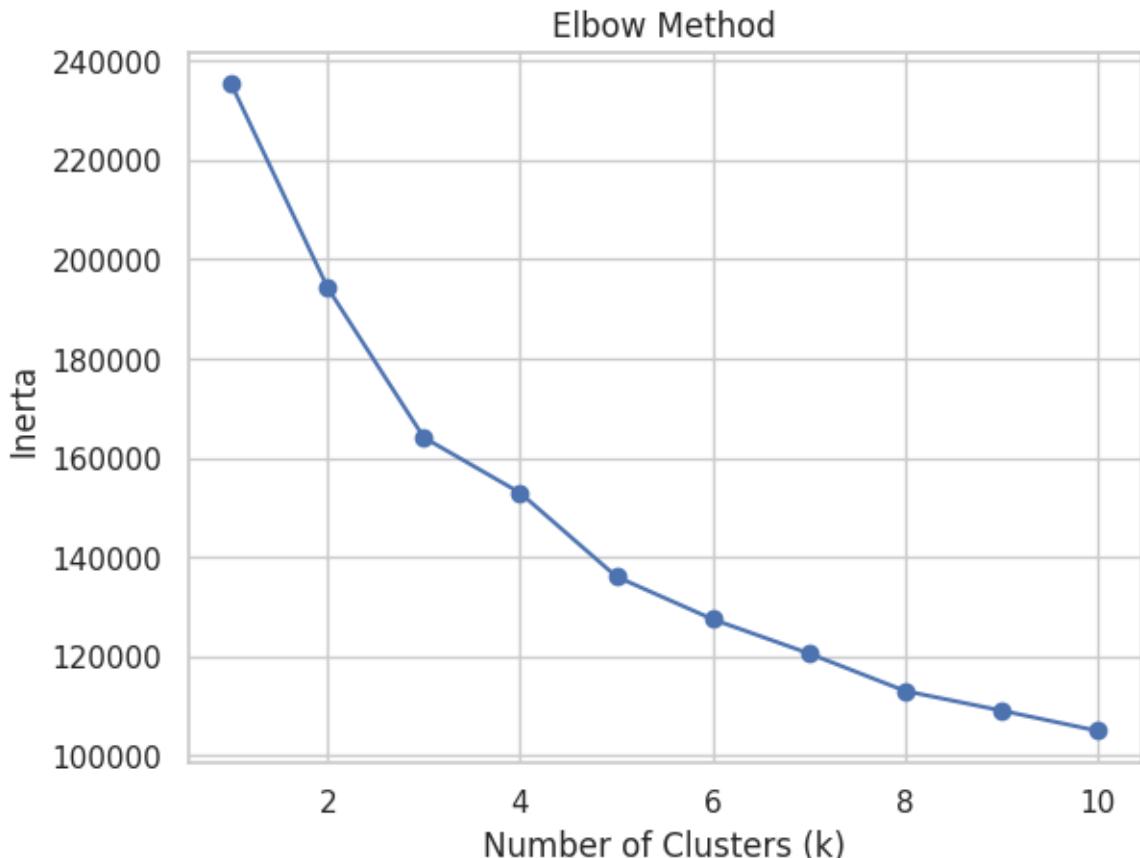
##### Elbow Method

- Use Elbow to select optimal clusters (k)
- Look for where decrease in inertia starts

In [61]:

```
wcss = []
for i in range(1, 11): # Try k values from 1 to 10
    kmeans = KMeans(n_clusters=i, random_state=42) # Initialize KMeans
    kmeans.fit(scaled_data_no_outliers) # Fit to your data
    wcss.append(kmeans.inertia_) # Append Within Cluster Sum of Squares (WCSS) aka inertia

plt.plot(range(1, 11), wcss, marker='o')
plt.title('Elbow Method')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia')
plt.show()
```



#### Conduct k-means

- Selected k as per elbow for feature selection.

```
K = 3
Choose K
In [62]:
```

```
# Choose k from Elbow graph where it bends up.
k = 3

# Perform KMeans clustering
kmeans_no_outliers = KMeans(n_clusters=k, random_state=42)
clusters_no_outliers = kmeans_no_outliers.fit_predict(scaled_data_no_outliers)

# Add cluster labels to the cleaned dataset
data_no_outliers['Cluster'] = clusters_no_outliers

# Display the dataset with the assigned clusters
data_no_outliers.head()
```

```
Out[62]:
```

				l	r	mi	num	ber	las	_of	t_r	e	rev	calcu	ai	roo	neigh	C
h	nei	ne	l	o	o	p	ni	ber	las	-	r	s_p	iew	lated	la	m_t	bourh	l
o	ghb	ig	a	n	o	mu				re	evi	er	er	host	bi	ype	ood_g	u
i	our	hb	hb	g	m	r				re	ew	v	mon	_list	li	_en	roup_	s
s			hoo	ou	i	i	m			re	vie	th	ings_	ty	cod	encod	e	r
t			d_g	rh	t	t	c	ni		re	ew	i	count	_3	ed	ed		
d			rou	oo	u	u	y	gh		re		e			65			
i			d	d	d	d	ts	ws				w						
d	p	d	e	e	e													

				h	a	av													
2	2	Man	Mi	.	e														
5	8	hat	dt	7	7	9	2	1		45	9-	05-	1	0.3	2	35	0	2	0
9	4	tan	ow	5	5	8	o	5			05-	8			5				
5	5		n	3	3	3	m				21								
				6	7	/													
				2	7	a													

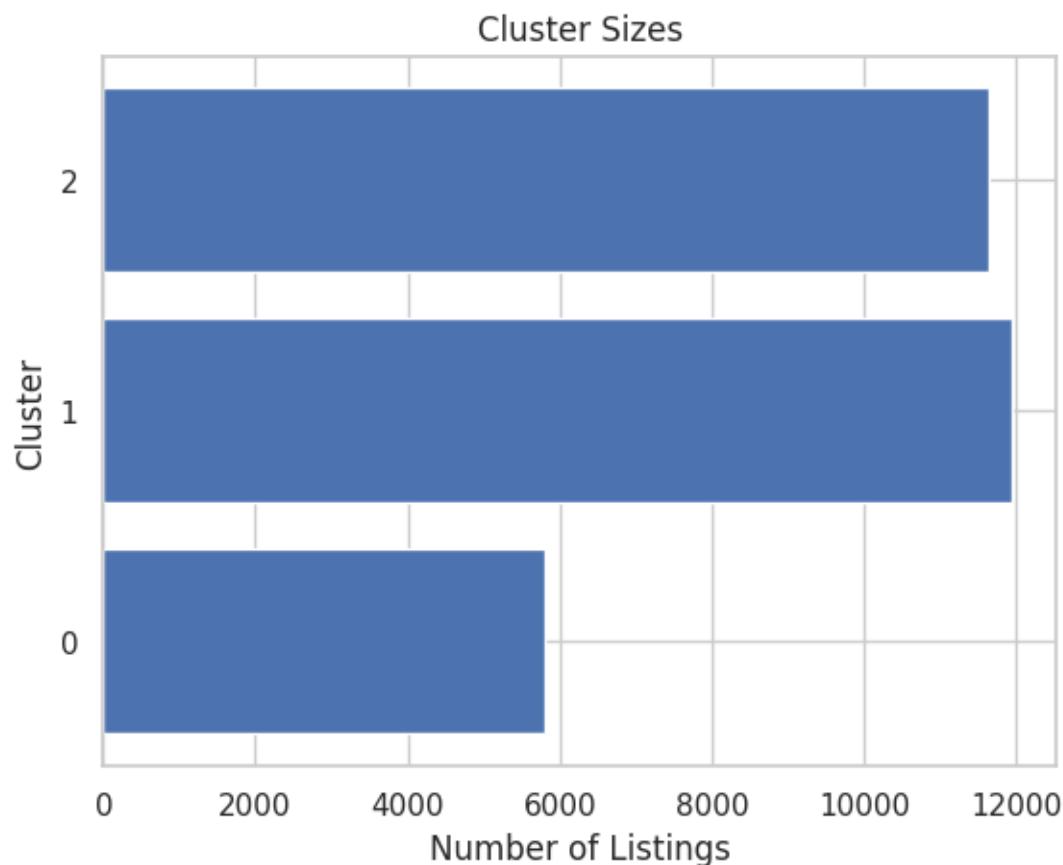
				E															
3	4	Man	Ha	.	4	7	i												
6	6	hat	rl	8	8	9	v												
4	3	tan	em	0	0	4	a	1	5	3	0	00:	0	0.0	1	36	1	2	2
7	2			9	9	1	t		0		0	0:	0			5			
				0	0	9	r				21								
				2	2	0	o				145								
							m				224								
											193								

				P															
5	7	Man	Ea	.	4	7	r												
0	1	hat	st	7	7	9	i												
2	9	tan	Ha	9	9	4		8	10	9									
2	2		rl	8	8	4		0			201								
			em	5	5	3	m	8-			7-								
				1	1	9	e	11-			09-								
						/	a	19			21								
							p				145								
							t				224								
											193								



```
plt.show()
```



**Calculate means to get insight into meaning of clusters.**  
Cluster 0: Popular, moderately priced and desirable

- Moderate price (\$117)
- Short minimum nights (2.5)
- High reviews (32) and per month (1.7)
- More than one listing (1.5)
- Available half of year (156)
- Mix room types but also entire homes
- Popular neighbourhood

Cluster 1: Desirable but expensive with limited availability

- Highest priced (\$166)
- Slightly longer minimum nights (3.3)
- Few reviews (6) and per month (.33)
- Just over one listing (1.1)
- Low availability (40)
- Mostly entire homes
- Popular neighbourhood

Cluster 2: Cheapest but less desirable with limited availability

- Cheapest (\$77)
- Short minimum nights (2.7)
- Fewest reviews (4) and per month (.27)
- Moderate listings (1.3)
- Moderately low availability (60)
- Mostly private or shared rooms
- Less popular or central neighbourhoods

In [65]:

```
cluster_means = data_no_outliers.groupby('Cluster').mean(numeric_only=True)
cluster_means
```

Out[65]:

	id	hos_t_i_d	la_ng_e	lo_ng_e	min_imu	numb_er_o_f_re	ha_s_ews	revi_ews_per	calculat_ed_host_listings	avai_lability	room_type	neighbo_urhood_group_encoded
C												
l												
u												
s												
t												
e												
r												
0	1.9 291 24e +07	7.3 343 94e +07	40 .7 28 48	- 73 .9 40	11 7. 37 40	2.4 764 40 04	32.8 2813 8 94	1. 00 00 8	1.73 4293 00 00	1.501808 2140 52	156. 2140 8672	0.52 1 1.66127
1	1.7 787 71e +07	4.8 747 74e +07	40 .7 29 23	- 73 .9 62	16 6. 34 03	3.3 046 97	5.78 5983 95	0. 76 0.33 73	40.4 7400 8002	1.110525 2	0.01 2	1.64523
2	1.8 964 97e +07	6.1 727 34e +07	40 .7 29 43	- 73 .9 44	77 .0 81 06	2.7 248 248 24	4.31 7550 40 66	0. 70 0.27 5226	59.5 7111 0	1.296620 8857	1.03 5	1.56776

#### Normalise means and visualise clusters

In [66]:

```
# features to visualise
features_to_visualize = ['price', 'minimum_nights', 'number_of_reviews', 'reviews_per_month',
                         'calculated_host_listings_count', 'availability_365',
                         'room_type_encoded', 'neighbourhood_group_encoded']

# Normalize the data using MinMaxScaler and store in 'normalised_cluster_means'
scaler = MinMaxScaler()
normalised_data = scaler.fit_transform(cluster_means[features_to_visualize])
normalised_cluster_means = pd.DataFrame(normalised_data, columns=features_to_visualize,
                                         index=cluster_means.index)

# Add the 'Cluster' column back to the normalized data
normalised_cluster_means['Cluster'] = cluster_means.index

# Explicitly cast all columns to float
normalised_cluster_means = normalised_cluster_means.astype(float)

# Add a small offset to zero values
normalised_cluster_means[normalised_cluster_means == 0] = 0.005 # Adjust the offset value as needed

# Set the width of the bars
bar_width = 0.1

# Set the positions of the bars on the x-axis (adjusted for more features)
r = np.arange(len(cluster_means.index))
positions = [r + i * bar_width for i in range(len(features_to_visualize))]

# Create the bar chart with adjusted figure size
plt.figure(figsize=(12, 6)) # Increased figure width for better visibility

# Color map for the bars
num_features = len(features_to_visualize)
cmap = plt.get_cmap('gist_rainbow', num_features)
colors = [cmap(i) for i in range(num_features)]
```

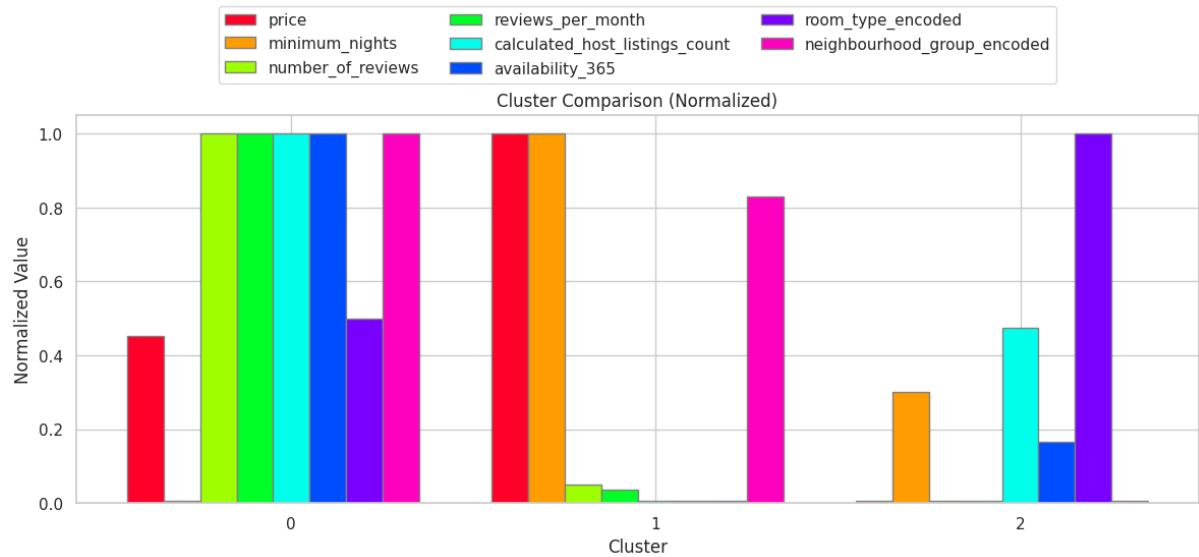
```

# Create bars for each feature
for i, feature in enumerate(features_to_visualize):
    plt.bar(positions[i], normalised_cluster_means[feature], color=colors[i], width=bar_width,
    edgecolor='grey', label=feature)

# Add labels, title, and legend
plt.xlabel('Cluster')
plt.ylabel('Normalized Value')
plt.title('Cluster Comparison (Normalized)')
plt.xticks([r + bar_width * (len(features_to_visualize) / 2) for r in
range(len(cluster_means.index))], cluster_means.index) # Center x-axis ticks
plt.legend(loc='upper center', bbox_to_anchor=(0.5, 1.3), ncol=3) # Adjust legend position

plt.tight_layout()
plt.show()

```



#### Visualise k-means using pairplot

- This takes a bit of time to run
- Creates pair-plots for the four features and five clusters

In [67]:

```

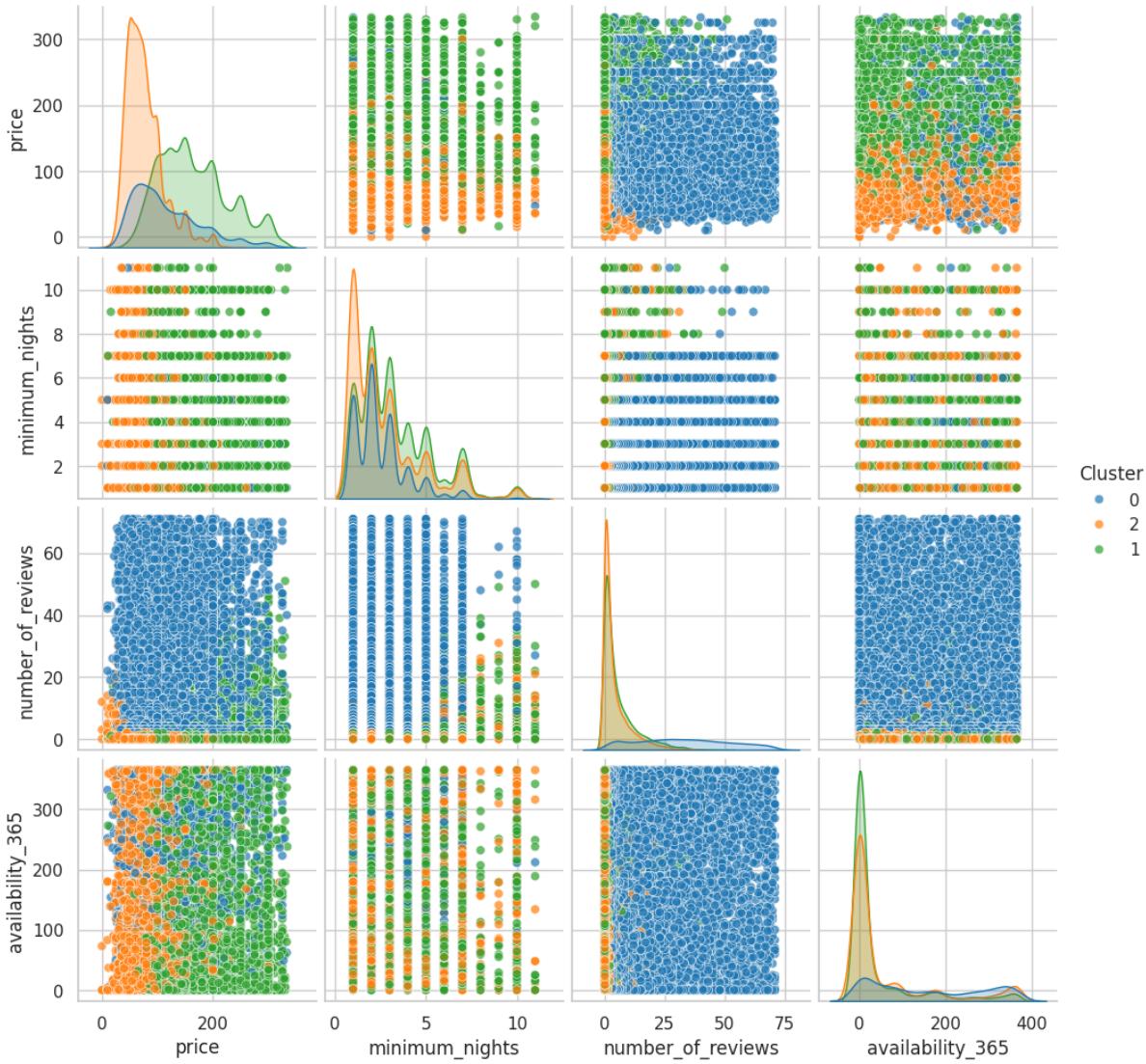
# Add cluster labels for visualization
data_no_outliers['Cluster'] = data_no_outliers['Cluster'].astype(str)

# Pairplot to visualize clusters based on selected features
sns.set(style="whitegrid")
pairplot = sns.pairplot(
    data_no_outliers,
    vars=['price', 'minimum_nights', 'number_of_reviews', 'availability_365'],
    hue='Cluster',
    palette='tab10',
    diag_kind='kde',
    plot_kws={'alpha': 0.7}
)

# Show the pairplot
plt.suptitle("KMeans Clustering Visualization", y=1.02)
plt.show()

```

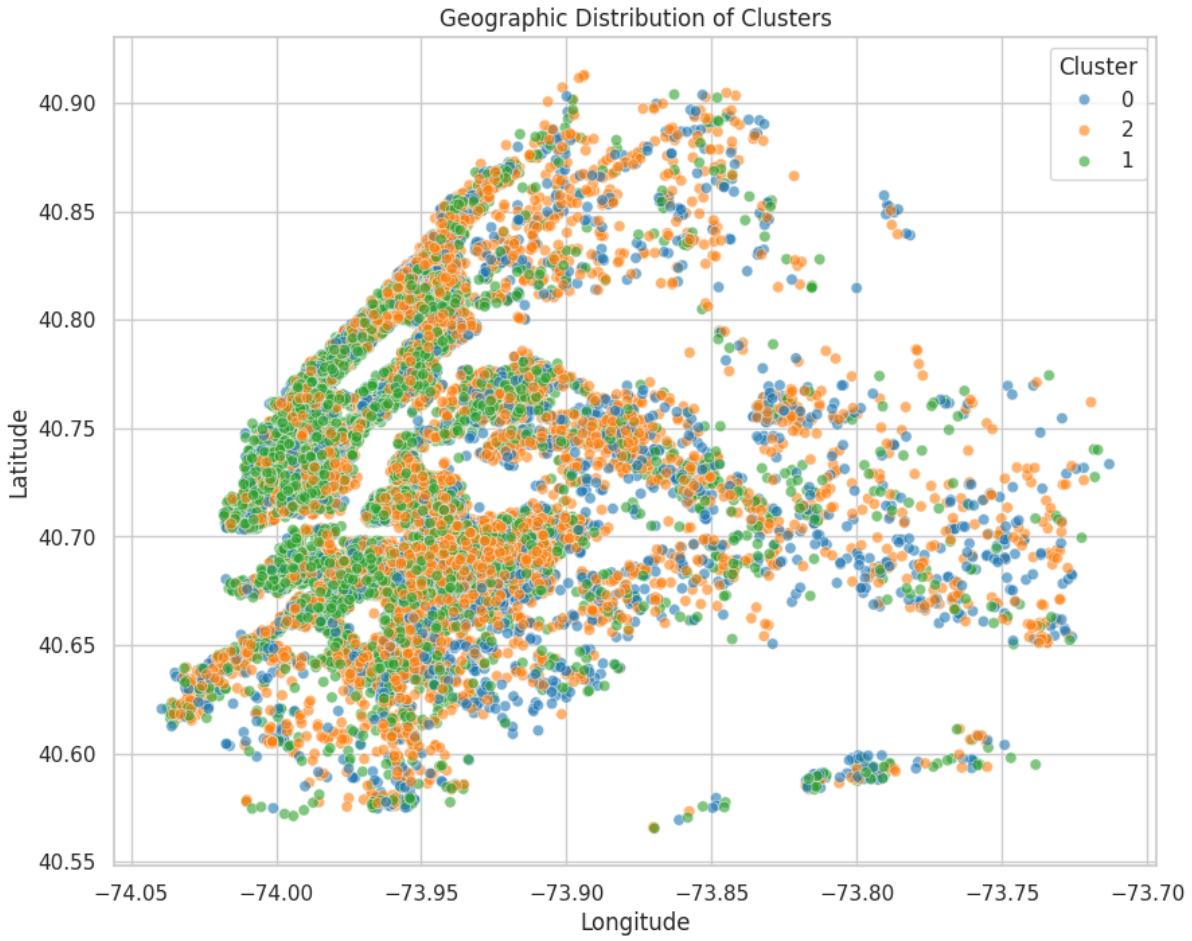
KMeans Clustering Visualization



#### Geographic distribution of clusters

In [68]:

```
# Geographic distribution of clusters
plt.figure(figsize=(10, 8))
sns.scatterplot(
    x=data_no_outliers['longitude'],
    y=data_no_outliers['latitude'],
    hue=data_no_outliers['Cluster'],
    palette='tab10',
    alpha=0.6
)
plt.title("Geographic Distribution of Clusters")
plt.xlabel("Longitude")
plt.ylabel("Latitude")
plt.legend(title="Cluster")
plt.show()
```



#### Geographic Visualisation

- Takes a bit of time to visualise.
- WARNING! WARNING! WARNING! If this has been run, can't check in to GitHub.

In [ ]:

```
# Create a GeoDataFrame for geographic visualization
geometry = [Point(xy) for xy in zip(data_no_outliers['longitude'],
data_no_outliers['latitude'])]
geo_df = gpd.GeoDataFrame(data_no_outliers, geometry=geometry)

# Base map creation using Folium
map_clusters = folium.Map(location=[data_no_outliers['latitude'].mean(),
data_no_outliers['longitude'].mean()], zoom_start=11)

# Add cluster points to the map
colors = ['blue', 'green', 'red']
for _, row in geo_df.iterrows():
    folium.CircleMarker(
        location=[row['latitude'], row['longitude']],
        radius=5,
        color=colors[int(row['Cluster'])]),
        fill=True,
        fill_opacity=0.6,
        popup=f"Cluster: {row['Cluster']}\nPrice: {row['price']}\nReviews: {row['number_of_reviews']}")
    .add_to(map_clusters)

# Display the map inline
map_clusters
```

#### REFERENCES

Ouleye, A. (2023) *Exploratory Data Analysis with Python Cookbook*. Pockt Publishing.