

steps

1.Importing all dependencies (It) 2.Loading datasets 3.Initial exploration 4.Data cleaning 5.Data Analysis

```
In [10]: #Importing all needed libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [11]: #Loading Dataset
data = pd.read_csv('datasets.csv')

In [16]: #Exploring the data
data.head()
```

```
Out[16]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	...	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_review
		Rental unit in Brooklyn	★5.0													
0	1312228+06	1	7150382	Walter	Brooklyn	Clinton Hill	40.683710	-73.964610	Private room	55.0	...	2012/15	0.03	1.0	0.0	
		bedroom														
		Rental unit in New York	★4.67													
1	4.527754e+07	2	51501835	Jennifer	Manhattan	Hell's Kitchen	40.766610	-73.988100	Entire home/apt	144.0	...	01/05/23	0.24	139.0	364.0	
		bedrooms														
		Rental unit in New York	★4.17													
2	9710000e+17	3	528871354	Joshua	Manhattan	Chelsea	40.750764	-73.994605	Entire home/apt	187.0	...	18/12/23	1.67	1.0	343.0	
		bedroom														
		Rental unit in New York	★4.64													
3	3.857863e+06	4	19902271	John And Catherine	Manhattan	Washington Heights	40.835660	-73.942500	Private room	120.0	...	17/09/23	1.38	2.0	363.0	
		bedroom														
		Condo in New York	★4.91													
4	4.089661e+07	5	61391963	Stay With Vibe	Manhattan	Murray Hill	40.761120	-73.978600	Entire home/apt	85.0	...	03/12/23	0.24	133.0	335.0	
		bed - 1..														

5 rows × 22 columns

data.tail()

```
In [8]: #Indicating Number of rows and columns
data.shape

Out[8]: (20770, 22)
```

```
In [9]: #Data information
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20770 entries, 0 to 20769
Data columns (total 22 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   id                   20770 non-null   float64
 1   name                 20770 non-null   object
 2   host_id              20770 non-null   int64
 3   host_name            20770 non-null   object
 4   neighbourhood_group  20770 non-null   object
 5   neighbourhood         20763 non-null   object
 6   latitude              20763 non-null   float64
 7   longitude             20763 non-null   float64
 8   room_type            20763 non-null   object
 9   price                20736 non-null   float64
10   minimum_nights       20763 non-null   float64
11   number_of_reviews    20763 non-null   float64
12   last_review          20763 non-null   object
13   reviews_per_month    20763 non-null   float64
14   calculated_host_listings_count  20763 non-null   float64
15   availability_365      20763 non-null   float64
16   number_of_reviews_ltn  20763 non-null   object
17   license              20770 non-null   object
18   rating               20770 non-null   object
19   bedrooms             20770 non-null   object
20   beds                 20767 non-null   int64
21   baths                20770 non-null   object
dtypes: float64(10), int64(2), object(10)
memory usage: 3.1+ MB

In [10]: #Statistical summary of data
data.describe()

Out[10]:
```

	id	host_id	latitude	longitude	price	minimum_nights	number_of_reviews	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds
count	2.077000e+04	2.077000e+04	20763.000000	20763.000000	20736.000000	20763.000000	20763.000000	20763.000000	20763.000000	20763.000000	20770.000000	20770.000000
mean	3.033558e+17	1.740949e+08	40.726821	-73.939179	187.714940	28.558493	42.610605	1.257589	18.666666	206.967967	10.849962	1.723592
std	3.901221e+17	1.725957e+08	0.080293	0.061403	1023.245124	33.532697	73.522401	1.904472	70.921443	135.077259	21.354876	1.211993
min	2.595000e+03	1.674000e+03	40.500314	-74.249840	10.000000	1.000000	1.000000	0.010000	1.000000	0.000000	0.000000	1.000000
25%	2.707800e+07	2.041184e+07	40.684159	-73.980765	80.000000	3.000000	4.000000	0.210000	1.000000	87.000000	1.000000	1.000000
50%	4.992852e+07	1.080990e+08	40.722890	-73.949597	125.000000	30.000000	14.000000	0.650000	2.000000	215.000000	3.000000	1.000000
75%	7.220000e+17	3.143979e+08	40.763106	-73.917475	199.000000	30.000000	49.000000	1.800000	5.000000	363.000000	15.000000	2.000000
max	1.050000e+18	5.504035e+08	40.911147	-73.713650	100000.000000	1250.000000	1865.000000	75.490000	713.000000	365.000000	1075.000000	42.000000

Data Cleaning

```
In [12]: #Checking Null values in each column
data.isnull().sum()

Out[12]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths
count	20770	20770	20770	20770	20770	20763	20763	20763	20763	20736	20763	20763	20763	20763	20763	20763	20767	20770	20770
dtype: object																			

```
In [20]: # dropping all missing values in rows
data.dropna(inplace=True)
data.isnull().sum()

Out[20]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths
count	20770	20770	20770	20770	20770	20763	20763	20763	20763	20736	20763	20763	20763	20763	20763	20763	20767	20770	20770
dtype: object																			

```
In [22]: # dealing with duplicates rows
data.duplicated().sum()

Out[22]: 12

In [24]: # deleting all duplicated rows
data.drop_duplicates(inplace=True)
data.duplicated().sum()

Out[24]: 0

In [26]: # type casting
# changing data types
data.dtypes

data['id'] = data['id'].astype(object)
data.dtypes

data['host_id'] = data['host_id'].astype(object)
data.dtypes

Out[26]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths
dtype: object	object	object	object	object	object	object	float64	float64	object	float64	float64	object	object	float64	float64	float64	int64	object	object

Applying Exploratory Data Analysis on the Dataset

Univariate Analysis

```
In [30]: # Identifying outliers in price
df = data[data['price'] < 1500]
sns.boxplot(data=df, y='price')

Out[30]: <Axes: ylabel='price'>
```



The box plot shows the distribution of prices. The y-axis is labeled 'price' and ranges from 0 to 1400. The x-axis is labeled 'price' and ranges from 0 to 1400. The plot shows a median price around 100, with a box from approximately 50 to 150. Whiskers extend from 0 to 1400. There are many outliers represented by small circles above the upper whisker, reaching up to 1500.

```
In [32]: #Price distribution
plt.figure(figsize=(8, 5))
sns.histplot(data=df, y='price', bins=100)
plt.title('Price Distribution')
plt.xlabel('price')
plt.ylabel('Frequency')
plt.show()

Out[32]:
```



The histogram shows the frequency of prices. The x-axis is labeled 'price' and ranges from 0 to 1400. The y-axis is labeled 'Frequency' and ranges from 0 to 1750. The distribution is right-skewed, with a peak frequency of approximately 1750 at a price of around 100. The frequency decreases as the price increases, with a long tail extending to 1400.

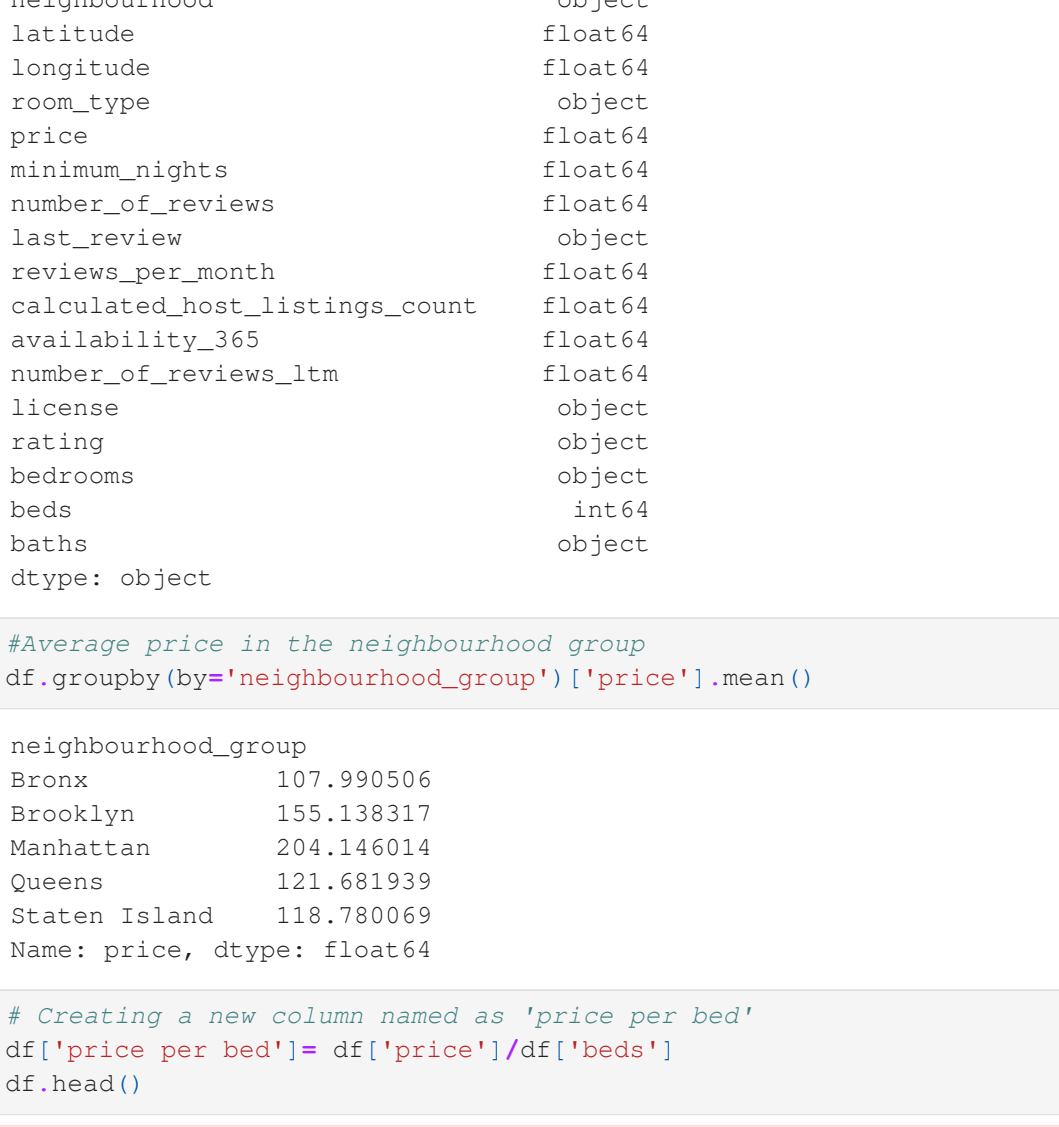
```
In [34]: df.dtypes

Out[34]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths
dtype: object	object	object	object	object	object	object	float64	float64	object	float64	float64	object	object	float64	float64	float64	int64	object	object

```
In [36]: #Price distribution
plt.figure(figsize=(6, 3))
sns.histplot(data=df, y='availability_365')
plt.title('availability_365 Distribution')
plt.xlabel('availability_365')
plt.ylabel('Frequency')
plt.show()

Out[36]:
```



The histogram shows the frequency of availability\_365. The x-axis is labeled 'availability\_365' and ranges from 0 to 350. The y-axis is labeled 'Frequency' and ranges from 0 to 5000. The distribution is right-skewed, with a peak frequency of approximately 5000 at an availability of around 350. The frequency decreases as the availability increases, with a long tail extending to 350.

```
In [38]: data.dtypes

Out[38]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths
dtype: object	object	object	object	object	object	object	float64	float64	object	float64	float64	object	object	float64	float64	float64	int64	object	object

```
In [40]: #Average price in the neighbourhood
df.groupby('neighbourhood_group')['price'].mean()

Out[40]:
```

neighbourhood_group	price
Bronx	107.995004
Brooklyn	155.138317
Manhattan	204.788493
Queens	121.481939
Staten Island	81.726101

```
In [42]: # Creating a new column named as 'price per bed'
df['price per bed'] = df['price']/df['beds']
df.head()

Out[42]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	...	reviews_per_month	calculated_host_listings_count	availability_365	number_of_review	price per bed
0	1312228.0	Rental unit in Brooklyn	★5.0		Brooklyn	Clinton Hill	40.683710	-73.964610	Private room	55.0	...	0.03	1.0	0.0	0.0	55.0
		bedroom														
1	4527757.0	Rental unit in New York	★4.67		Manhattan	Hell's Kitchen	40.766610	-73.988100	Entire home/apt	144.0	...	0.24	139.0	364.0	2	72.0
		bedrooms														
2	97100000.00000000.0	Rental unit in New York	★4.17		Manhattan	Chelsea	40.750764	-73.994605	Entire home/apt	187.0	...	1.67	1.0	343.0	6	31.166666666666664
		bedroom														
3	3857863.0	Rental unit in New York	★4.64		Manhattan	Washington Heights	40.835660	-73.942500	Private room	120.0	...	1.38	2.0	363.0	12	10.0
		bedroom														
4	40896611.0	Condo in New York	★4.91		Manhattan	Murray Hill	40.761120	-73.978600	Entire home/apt	85.0	...	0.24	133.0	335.0	3	28.333333333333332
		bed - 1..														

```
In [44]: # Average price per bed
df.groupby('neighbourhood_group')['price per bed'].mean()

Out[44]:
```

neighbourhood_group	price per bed
Bronx	74.713639
Brooklyn	89.788493
Manhattan	139.709057
Queens	76.236210
Staten Island	57.126101

Bivariate Analysis

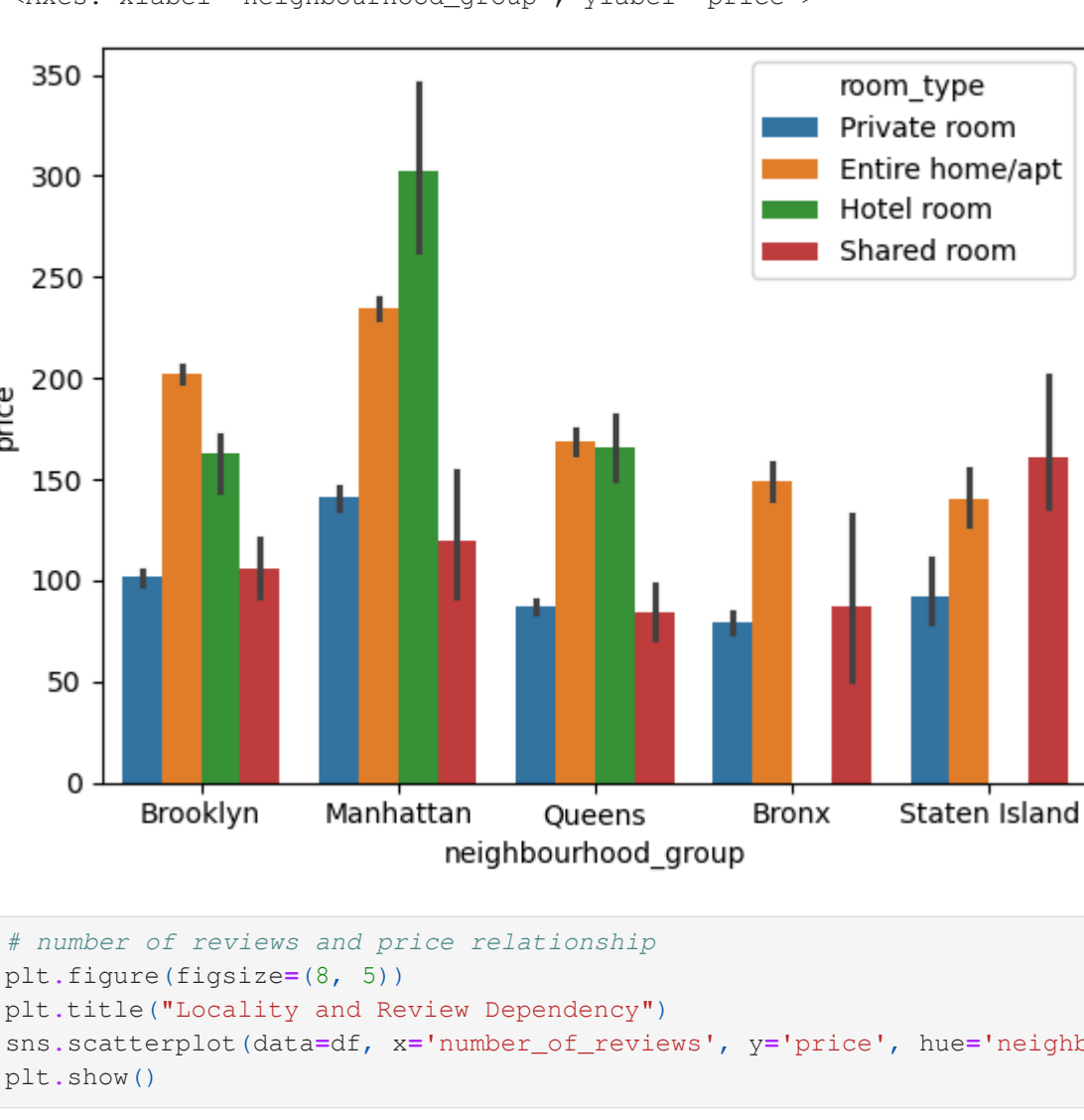
```
In [47]: df.columns

Out[47]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths	price per bed
dtype: object	object	object	object	object	object	object	float64	float64	object	float64	float64	object	object	float64	float64	float64	int64	object	object	float64

```
In [49]: # price dependency on neighbourhood
sns.pairplot(data=df, y='price', hue='neighbourhood_group', y_vars='price', hue='room_type')

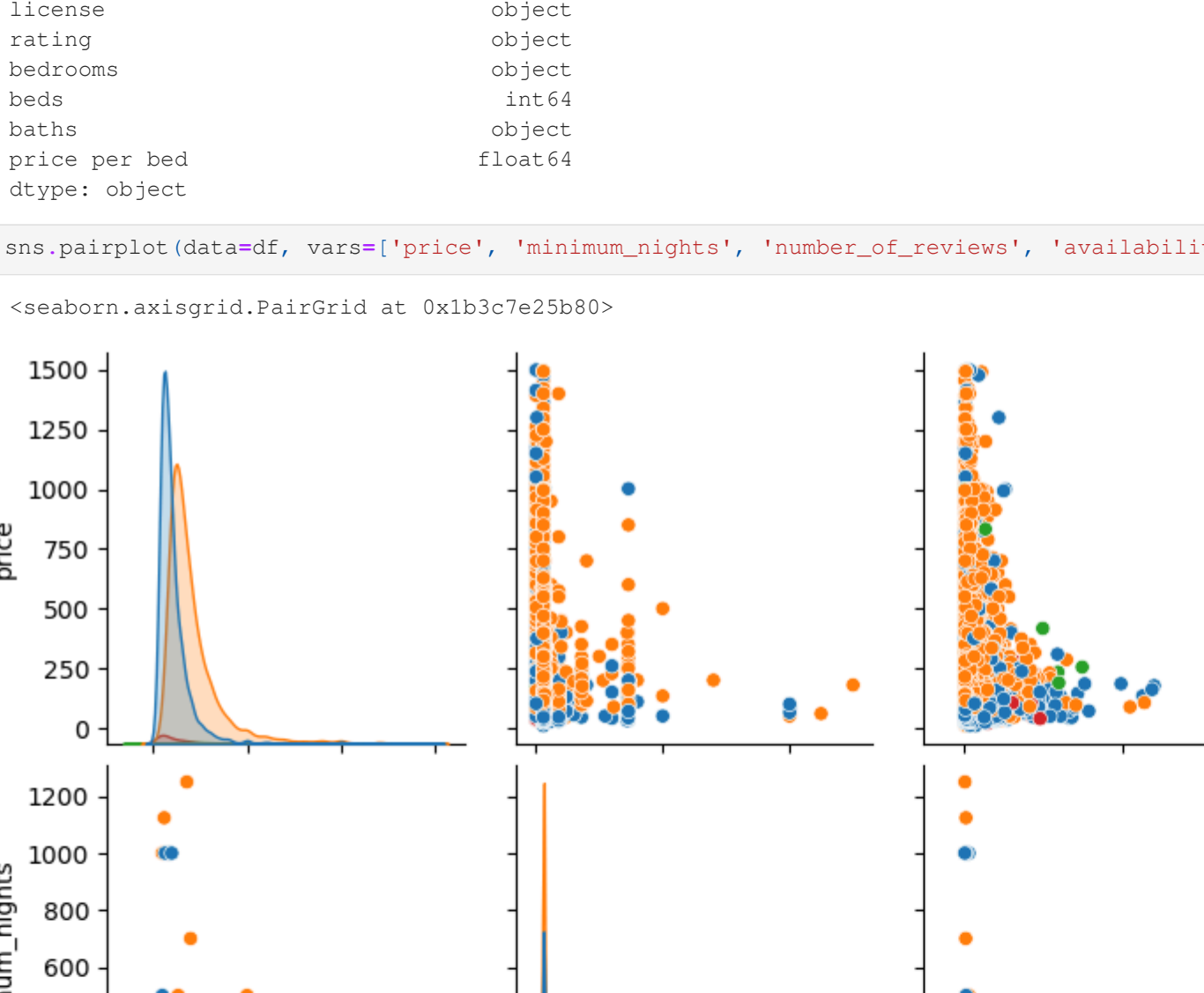
Out[49]: <Axes: xlabel='neighbourhood_group', ylabel='price'>
```



The pair plot shows the relationship between price and neighbourhood\_group. The x-axis is labeled 'neighbourhood\_group' and has categories: Brooklyn, Manhattan, Queens, Bronx, and Staten Island. The y-axis is labeled 'price' and ranges from 0 to 350. The plot shows that price generally increases with the neighbourhood\_group, with Manhattan having the highest prices and Staten Island having the lowest. The plot also shows the relationship between price and room\_type, with different room types having different price distributions.

```
In [53]: # number of reviews and price relationship
plt.figure(figsize=(8, 5))
plt.title('Locality and Review Dependency')
sns.scatterplot(data=df, x='number_of_reviews', y='price', hue='neighbourhood_group')
plt.show()

Out[53]:
```



The scatter plot shows the relationship between price and number\_of\_reviews. The x-axis is labeled 'number\_of\_reviews' and ranges from 0 to 1750. The y-axis is labeled 'price' and ranges from 0 to 1400. The plot shows that price generally increases with the number\_of\_reviews, with a positive correlation. The plot also shows the relationship between price and neighbourhood\_group, with different neighbourhoods having different price distributions.

```
In [55]: df.dtypes

Out[55]:
```

	id	name	host_id	host_name	neighbourhood_group	neighbourhood	latitude	longitude	room_type	price	minimum_nights	number_of_reviews	last_review	reviews_per_month	calculated_host_listings_count	availability_365	number_of_reviews_ltn	beds	baths	price per bed
dtype: object	object	object	object	object	object	object	float64	float64	object	float64	float64	object	object	float64	float64	float64	int64	object	object	float64

```
In [57]: sns.pairplot(data=df, vars=['price', 'minimum_nights', 'number_of_reviews', 'availability_365', 'room_type'])

Out[57]:
```



The pair plot shows the relationships between price, minimum\_nights, number\_of\_reviews, availability\_365, and room\_type. The plot shows that price generally increases with the number\_of\_reviews, with a positive correlation. The plot also shows the relationship between price and room\_type, with different room types having different price distributions. The plot also shows the relationship between minimum\_nights and number\_of\_reviews, with a negative correlation. The plot also shows the relationship between availability\_365 and number\_of\_reviews, with a positive correlation.

```
In [59]: #Geographical Distribution of Airbnb Listing
plt.figure(figsize=(10, 7))
sns.scatterplot(data=df, x='longitude', y='latitude', hue='room_type')
plt.title('Geographical Distribution of Airbnb Listing')
plt.show()

Out[59]:
```



The scatter plot shows the geographical distribution of Airbnb listings. The x-axis is labeled 'longitude' and ranges from -74.2 to -73.7. The y-axis is labeled 'latitude' and ranges from 40.5 to 40.9. The plot shows that the listings are concentrated in the New York City area, with a high density of listings in the Manhattan area. The plot also shows the relationship between longitude and latitude, with a positive correlation.

```
In [61]: #Heatmap - correlation of one variable with others for numerical column
corr = df[['latitude', 'longitude', 'price', 'minimum_nights', 'number_of_reviews', 'reviews_per_month', 'availability_365', 'beds']].corr()
corr

Out[61]: <Axes: >
```



The heatmap shows the correlation matrix for the numerical variables. The x-axis and y-axis are labeled with the variable names: latitude, longitude, price, minimum\_nights, number\_of\_reviews, reviews\_per\_month, availability\_365, and beds. The plot shows that the variables are highly correlated, with a strong positive correlation between price and number\_of\_reviews, and a strong negative correlation between minimum\_nights and number\_of\_reviews.



