

# house-rent-eda

February 3, 2024

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[2]: df = pd.read_csv('House_Rent_Dataset.csv')
```

```
[3]: df.head()
```

```
[3]:
```

	Posted On	BHK	Rent	Size	Floor	Area Type	\
0	2022-05-18	2	10000	1100	Ground out of 2	Super Area	
1	2022-05-13	2	20000	800	1 out of 3	Super Area	
2	2022-05-16	2	17000	1000	1 out of 3	Super Area	
3	2022-07-04	2	10000	800	1 out of 2	Super Area	
4	2022-05-09	2	7500	850	1 out of 2	Carpet Area	

	Area Locality	City	Furnishing Status	Tenant Preferred	\
0	Bandel	Kolkata	Unfurnished	Bachelors/Family	
1	Phool Bagan, Kankurgachi	Kolkata	Semi-Furnished	Bachelors/Family	
2	Salt Lake City Sector 2	Kolkata	Semi-Furnished	Bachelors/Family	
3	Dumdum Park	Kolkata	Unfurnished	Bachelors/Family	
4	South Dum Dum	Kolkata	Unfurnished	Bachelors	

	Bathroom	Point of Contact
0	2	Contact Owner
1	1	Contact Owner
2	1	Contact Owner
3	1	Contact Owner
4	1	Contact Owner

```
[4]: df.columns
```

```
[4]: Index(['Posted On', 'BHK', 'Rent', 'Size', 'Floor', 'Area Type',
          'Area Locality', 'City', 'Furnishing Status', 'Tenant Preferred',
          'Bathroom', 'Point of Contact'],
          dtype='object')
```

```
[4]: df.isnull().sum()
```

```
[4]: Posted On      0
      BHK           0
      Rent          0
      Size          0
      Floor         0
      Area Type     0
      Area Locality 0
      City          0
      Furnishing Status 0
      Tenant Preferred 0
      Bathroom      0
      Point of Contact 0
      dtype: int64
```

```
[5]: df.describe()
```

```
[5]:
```

	BHK	Rent	Size	Bathroom
count	4746.000000	4.746000e+03	4746.000000	4746.000000
mean	2.083860	3.499345e+04	967.490729	1.965866
std	0.832256	7.810641e+04	634.202328	0.884532
min	1.000000	1.200000e+03	10.000000	1.000000
25%	2.000000	1.000000e+04	550.000000	1.000000
50%	2.000000	1.600000e+04	850.000000	2.000000
75%	3.000000	3.300000e+04	1200.000000	2.000000
max	6.000000	3.500000e+06	8000.000000	10.000000

```
[6]: df.corr()
```

C:\Users\Admin\AppData\Local\Temp\ipykernel\_7396\1134722465.py:1: FutureWarning:  
The default value of numeric\_only in DataFrame.corr is deprecated. In a future  
version, it will default to False. Select only valid columns or specify the  
value of numeric\_only to silence this warning.

```
df.corr()
```

```
[6]:
```

	BHK	Rent	Size	Bathroom
BHK	1.000000	0.369718	0.716145	0.794885
Rent	0.369718	1.000000	0.413551	0.441215
Size	0.716145	0.413551	1.000000	0.740703
Bathroom	0.794885	0.441215	0.740703	1.000000

```
[7]: df.rename(columns={'Posted On': 'Date'}, inplace=True)
```

```
[8]: df.head(1)
```

```
[8]:
```

	Date	BHK	Rent	Size	Floor	Area Type	Area Locality \
0	2022-05-18	2	10000	1100	Ground out of 2	Super Area	Bandel

	City	Furnishing Status	Tenant Preferred	Bathroom	Point of Contact
0	Kolkata	Unfurnished	Bachelors/Family	2	Contact Owner

```
[9]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4746 entries, 0 to 4745
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Date                  4746 non-null  object
1   BHK                   4746 non-null  int64
2   Rent                  4746 non-null  int64
3   Size                  4746 non-null  int64
4   Floor                 4746 non-null  object
5   Area Type             4746 non-null  object
6   Area Locality         4746 non-null  object
7   City                  4746 non-null  object
8   Furnishing Status     4746 non-null  object
9   Tenant Preferred      4746 non-null  object
10  Bathroom              4746 non-null  int64
11  Point of Contact      4746 non-null  object
dtypes: int64(4), object(8)
memory usage: 445.1+ KB
```

```
[10]: df['Date'] = pd.to_datetime(df['Date'])
```

```
[11]: df['Date'].dtype
```

```
[11]: dtype('<M8[ns]')
```

```
[12]: # create a separate column of day and month
df['Day'] = df['Date'].dt.day
df['Month'] = df['Date'].dt.month
```

```
[13]: df.head(3)
```

```
[13]:
```

	Date	BHK	Rent	Size	Floor	Area Type	Area Locality \
0	2022-05-18	2	10000	1100	Ground out of 2	Super Area	Bandel
1	2022-05-13	2	20000	800	1 out of 3	Super Area	Kolkata
2	2022-05-16	2	17000	1000	1 out of 3	Super Area	Unfurnished

	City	Furnishing Status	Tenant Preferred
0	Kolkata	Unfurnished	Bachelors/Family

1	Phool Bagan, Kankurgachi	Kolkata	Semi-Furnished	Bachelors/Family
2	Salt Lake City Sector 2	Kolkata	Semi-Furnished	Bachelors/Family

	Bathroom	Point of Contact	Day	Month
0	2	Contact Owner	18	5
1	1	Contact Owner	13	5
2	1	Contact Owner	16	5

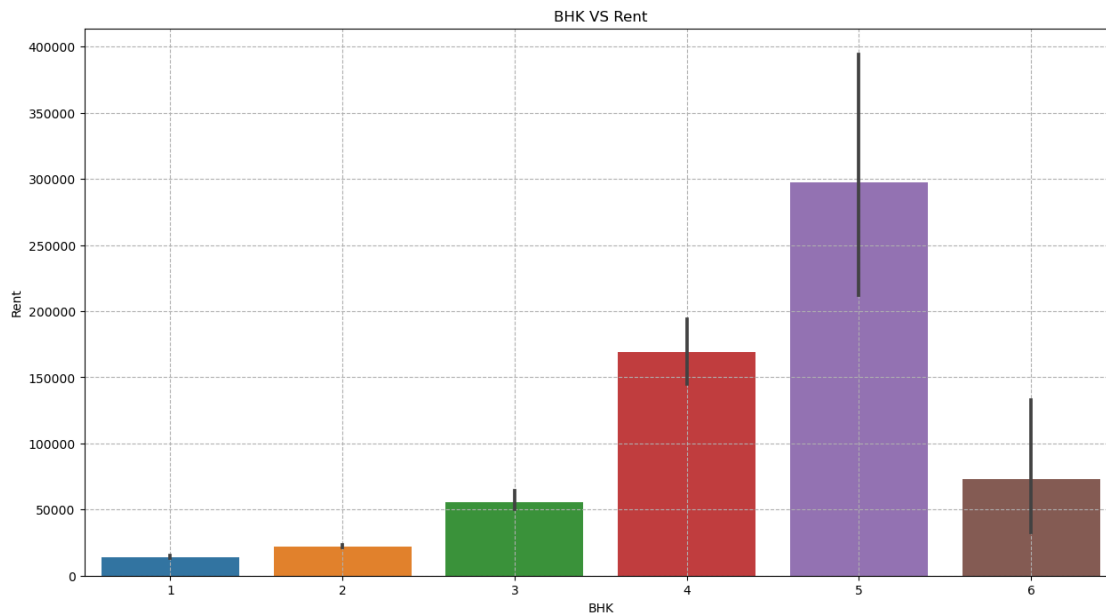
```
[14]: df.drop('Date', axis=1, inplace=True)
```

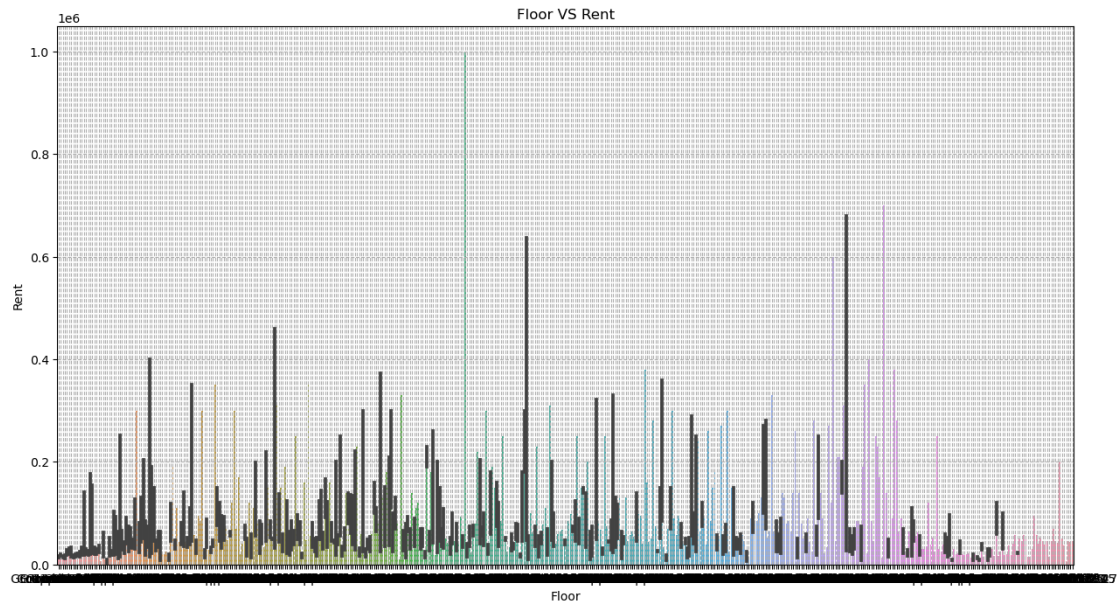
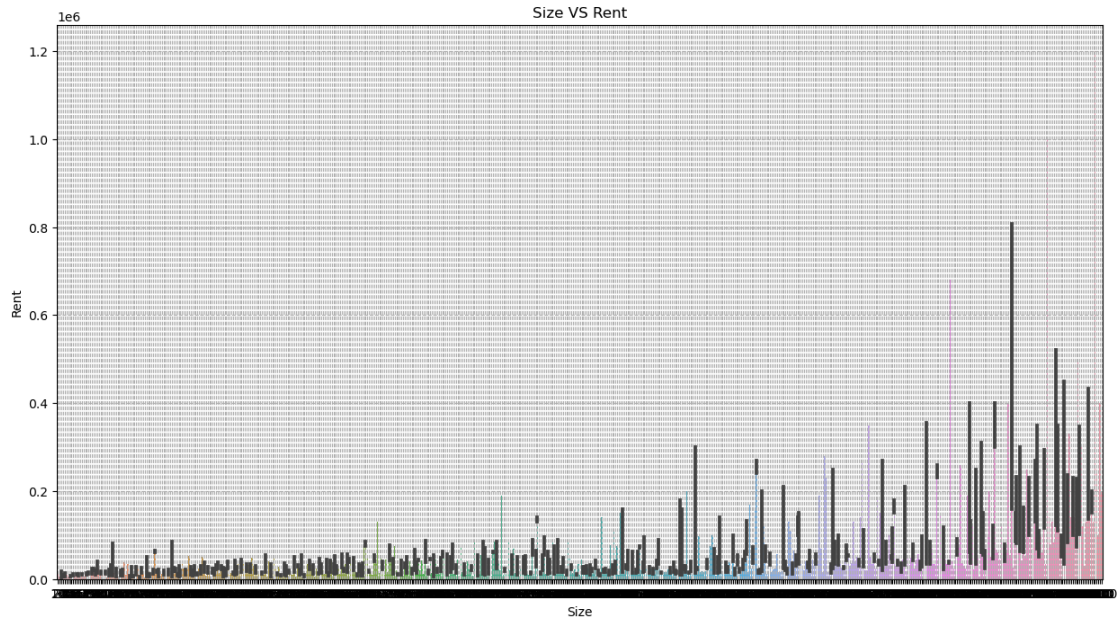
```
[15]: df.head(1)
```

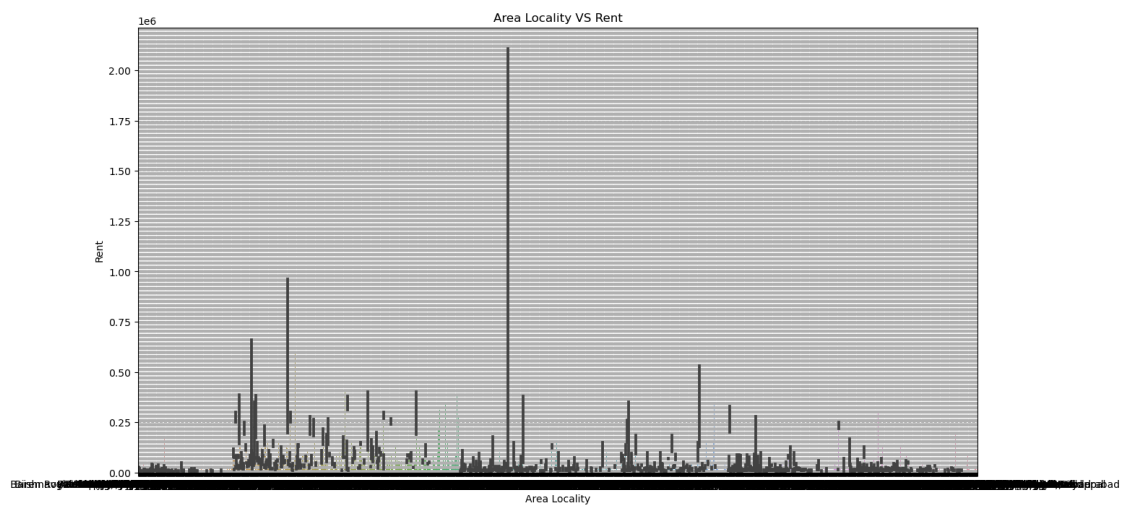
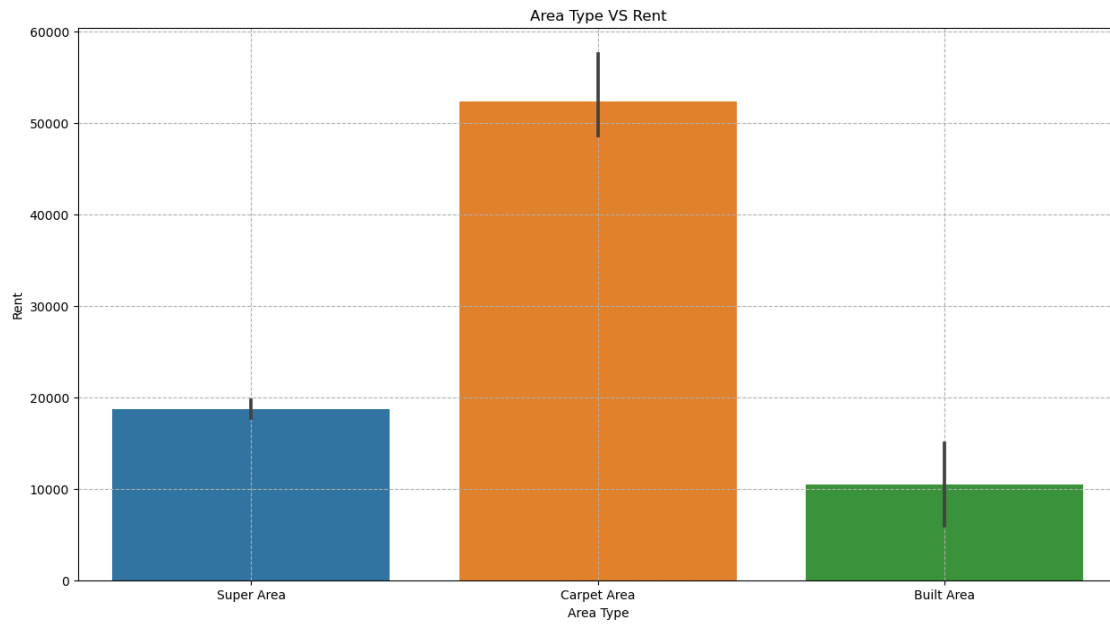
```
[15]:   BHK   Rent  Size      Floor  Area Type Area Locality   City \
0    2  10000  1100  Ground out of 2  Super Area      Bandel  Kolkata
```

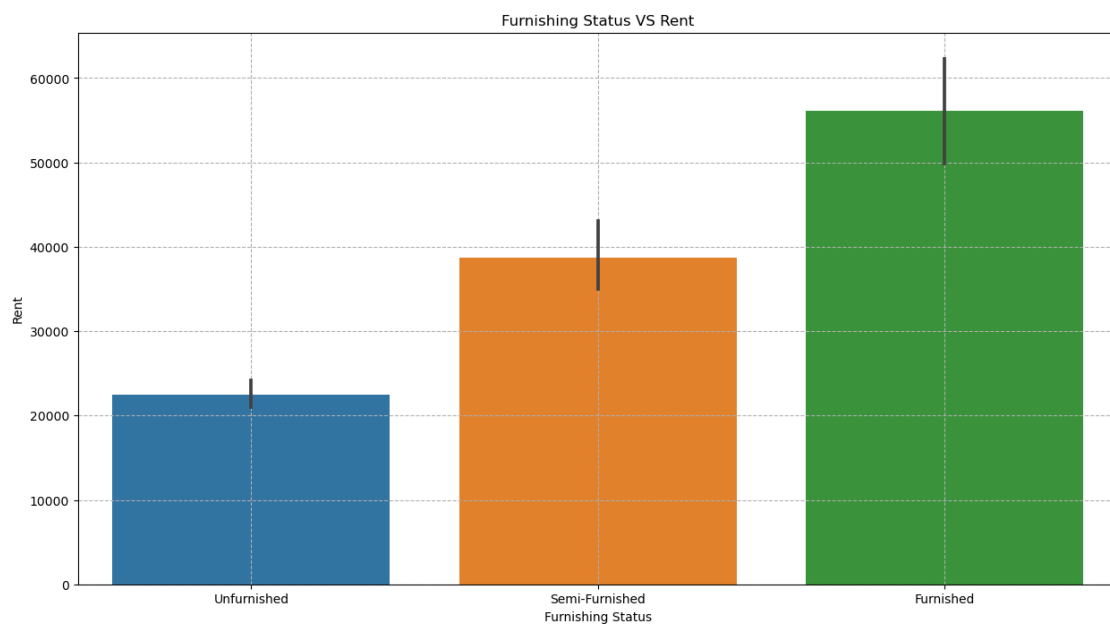
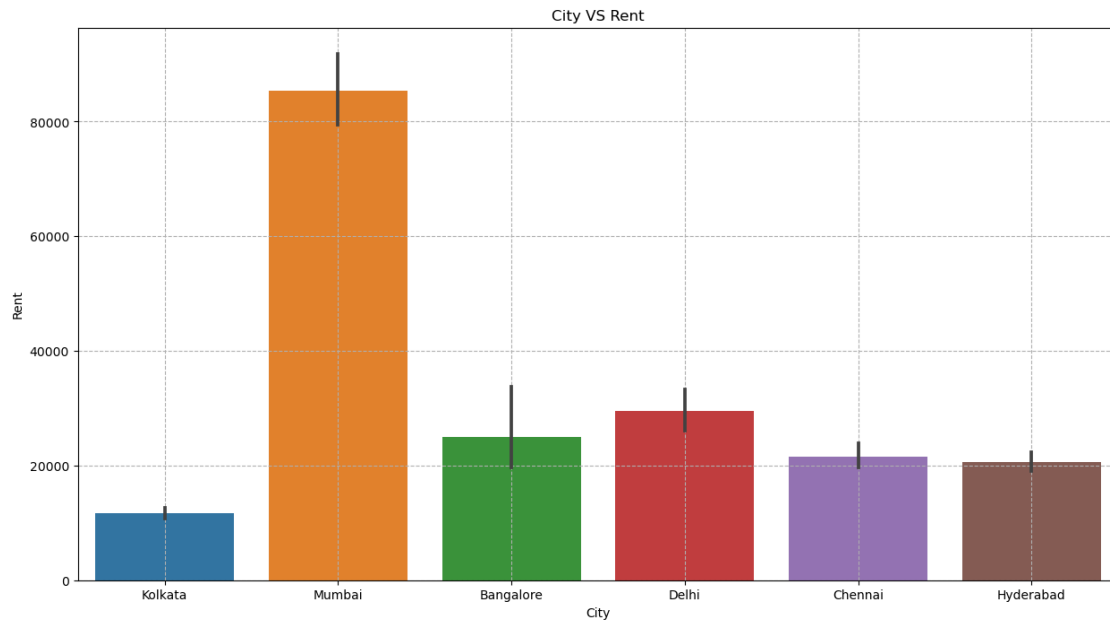
	Furnishing	Status	Tenant Preferred	Bathroom	Point of Contact	Day	Month
0	Unfurnished	Bachelors/Family		2	Contact Owner	18	5

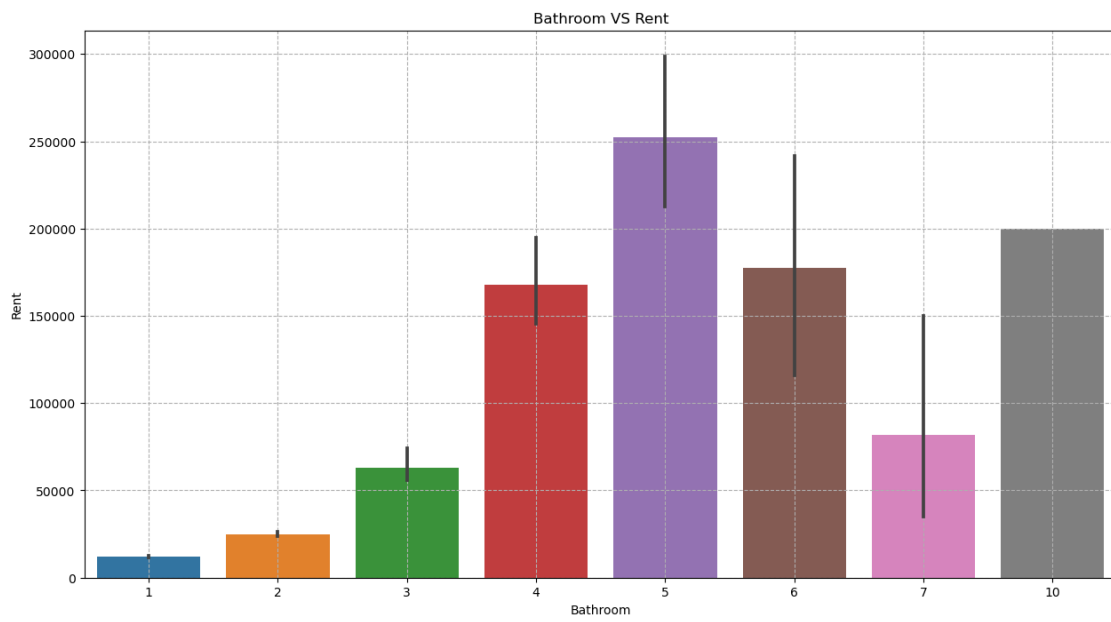
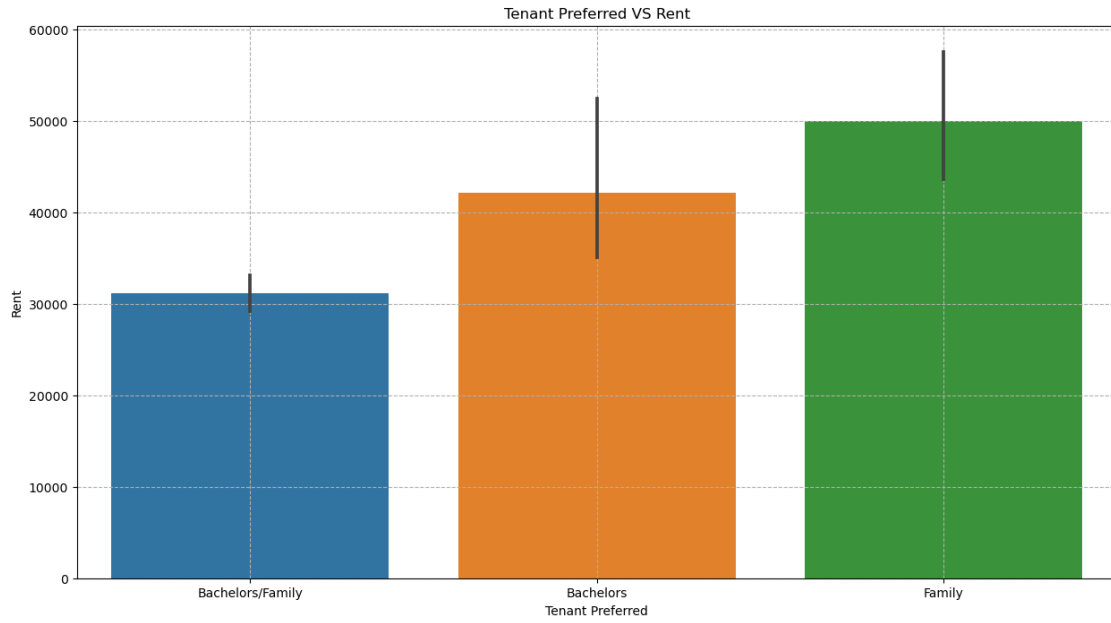
```
[19]: target_column = 'Rent'
sub_columns = [p for p in df.columns if p != target_column]
for sub_column in sub_columns:
    plt.figure(figsize=(15,8))
    sns.barplot(x = sub_column, y = target_column, data= df)
    plt.title(f'{sub_column} VS {target_column}')
    plt.grid(linestyle='--')
    plt.show()
```



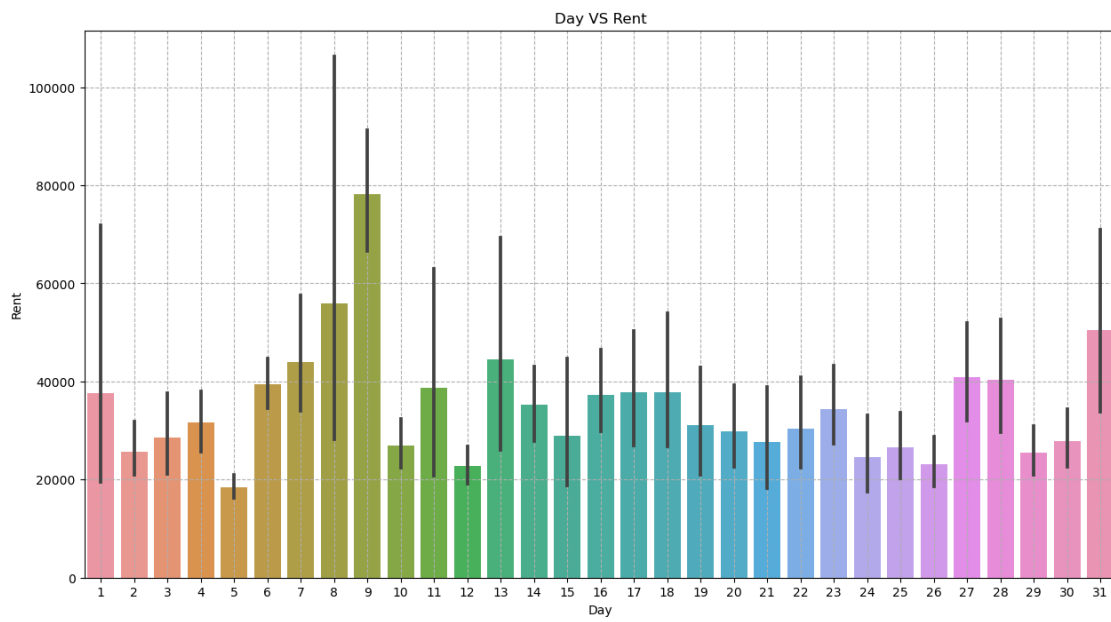
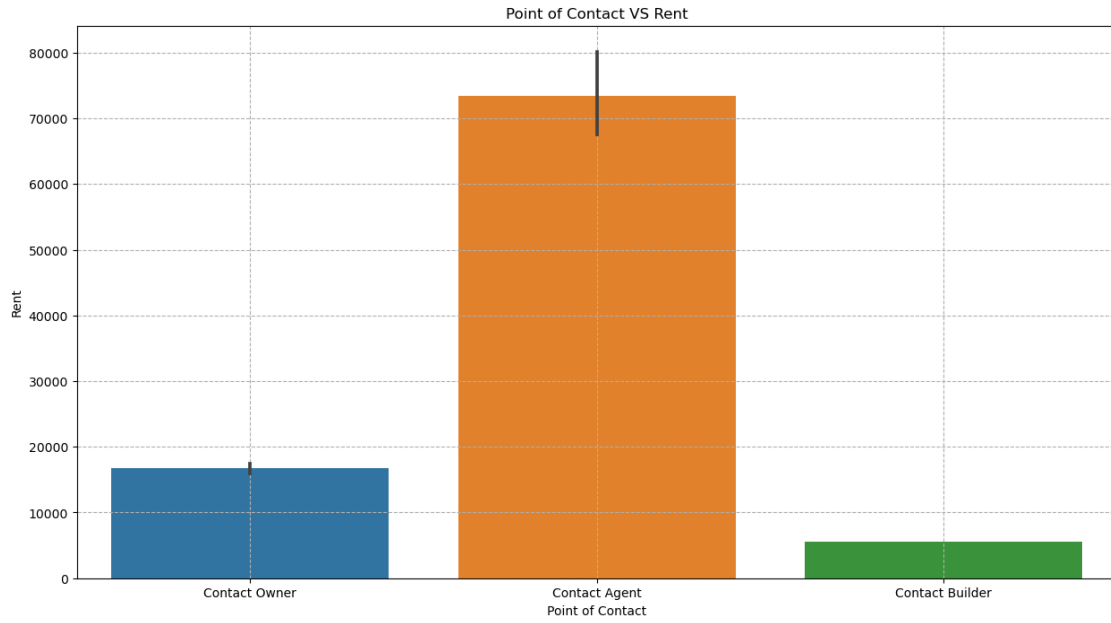


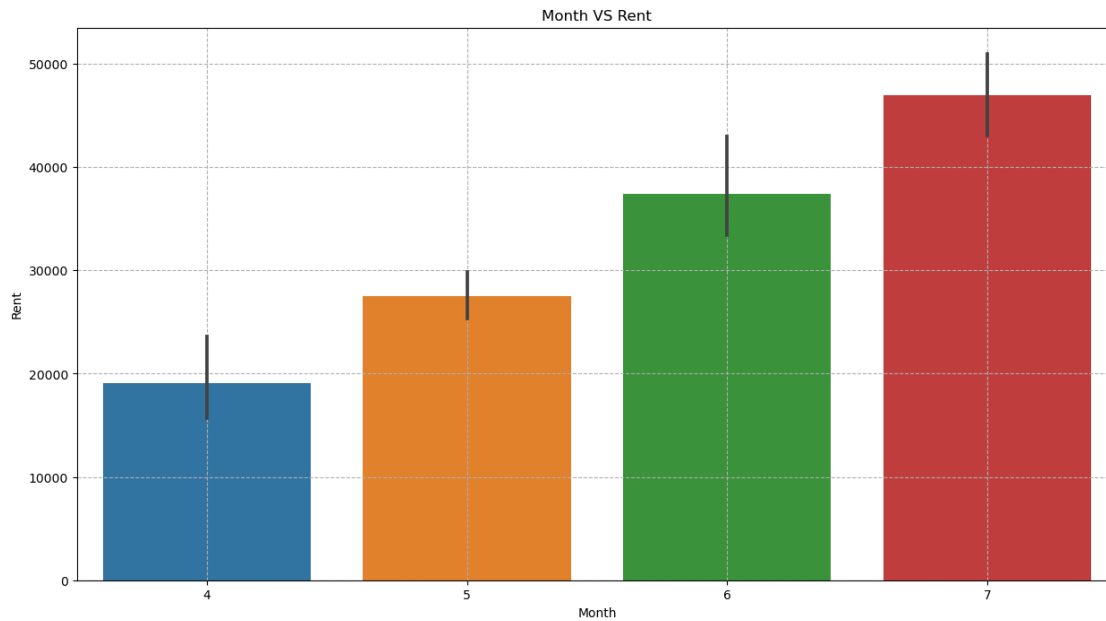




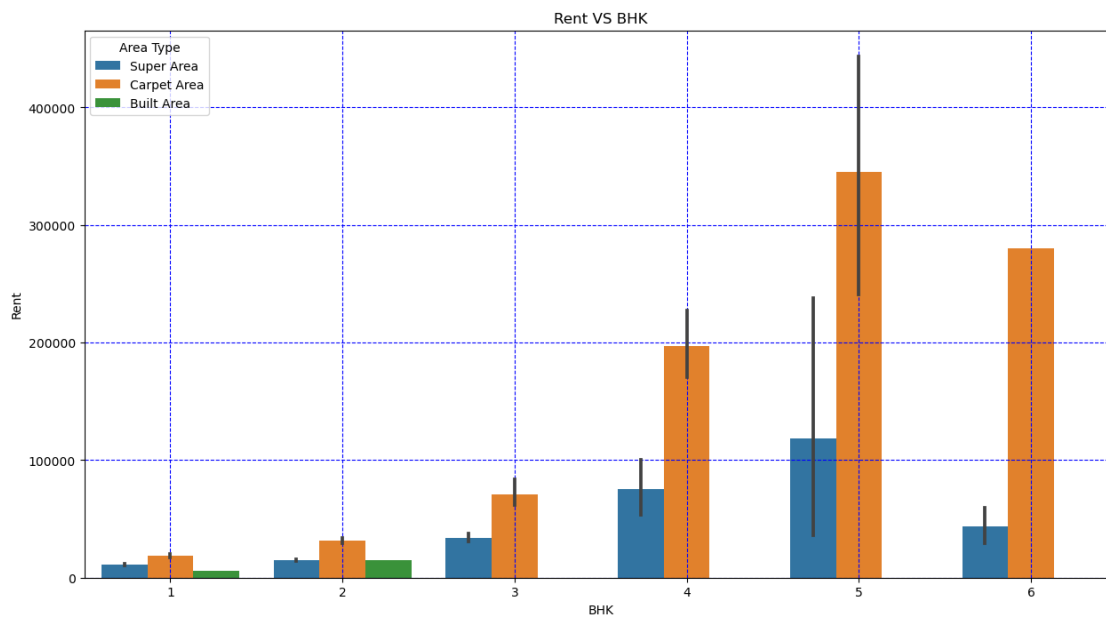


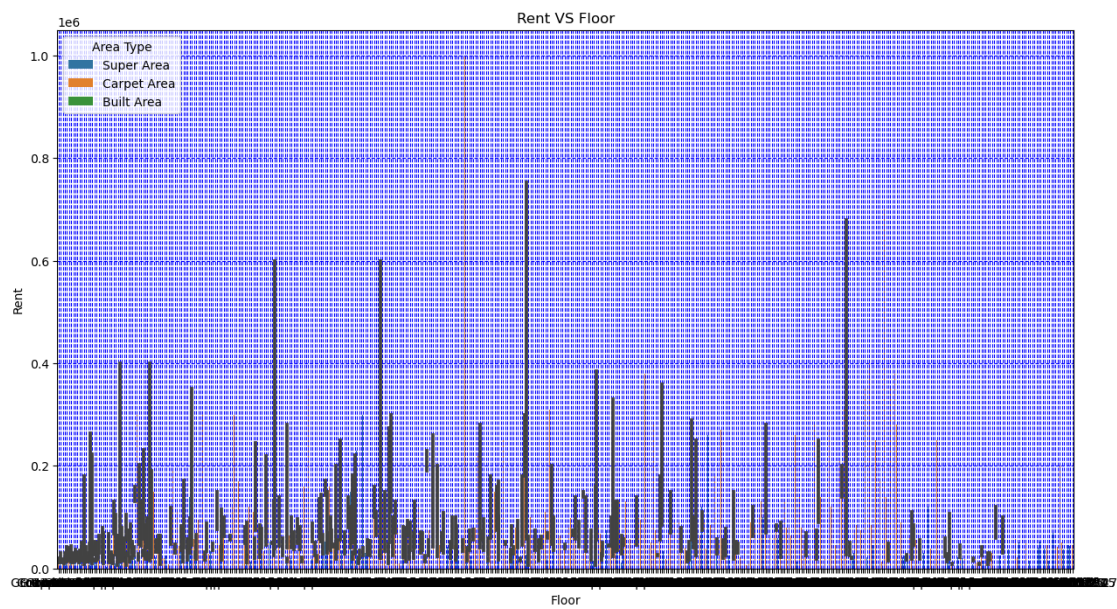
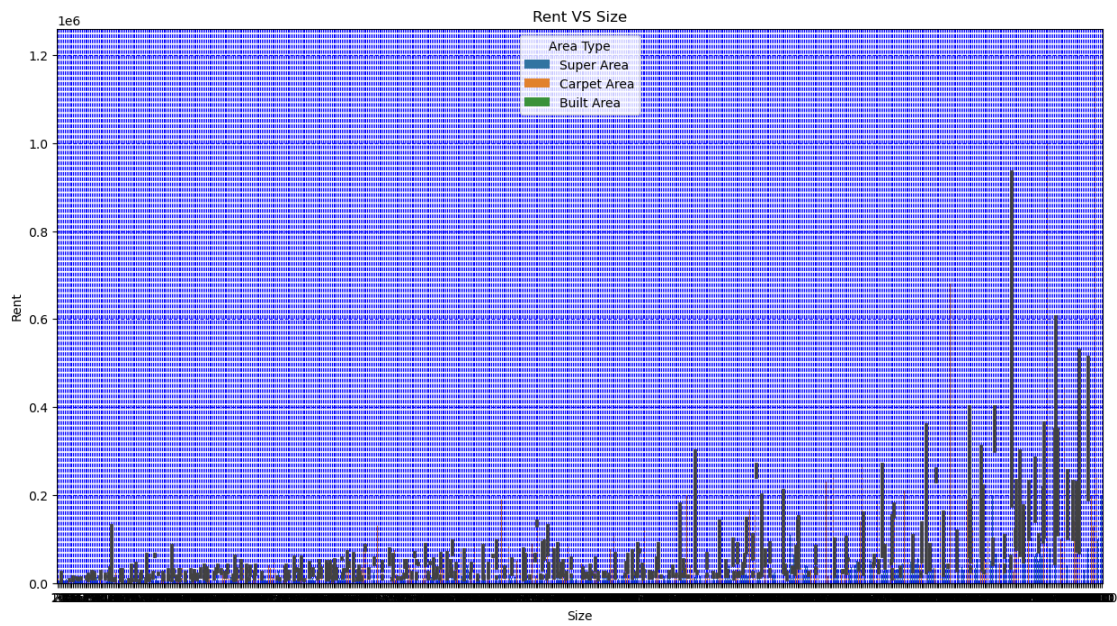


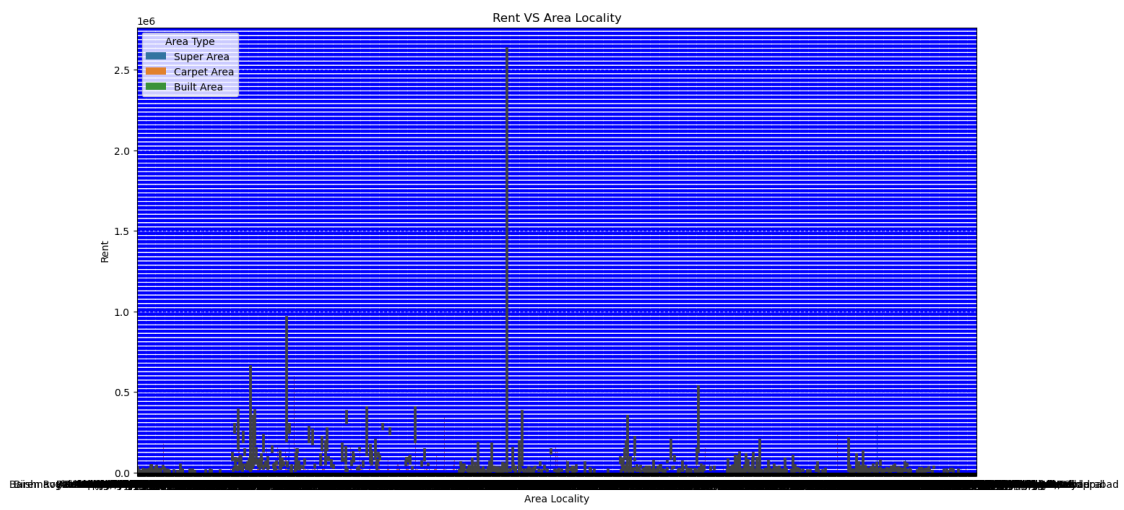
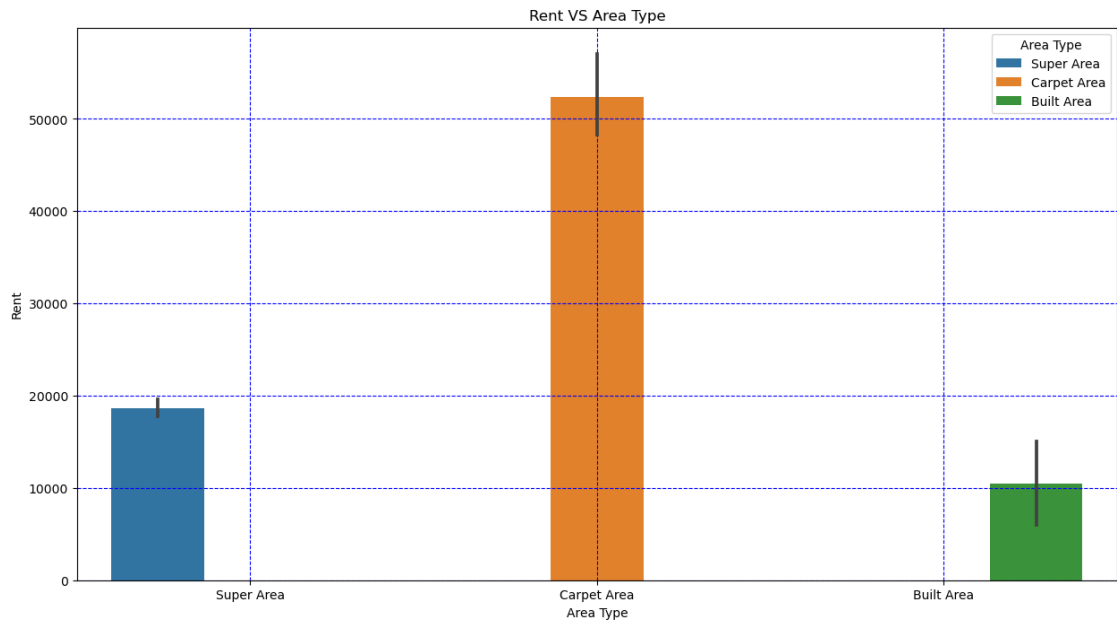


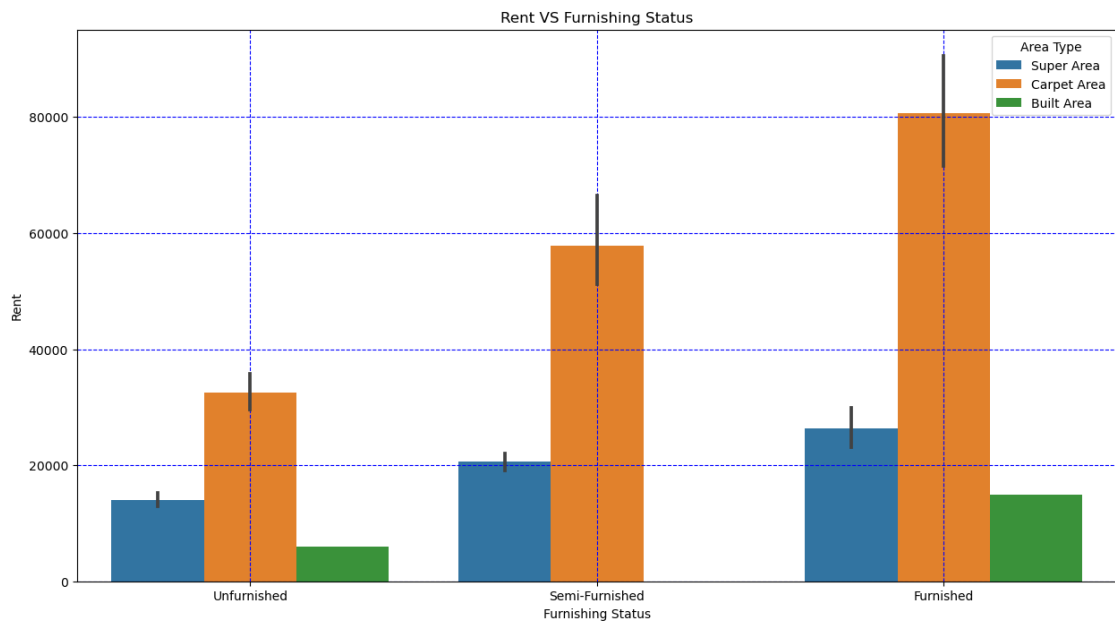
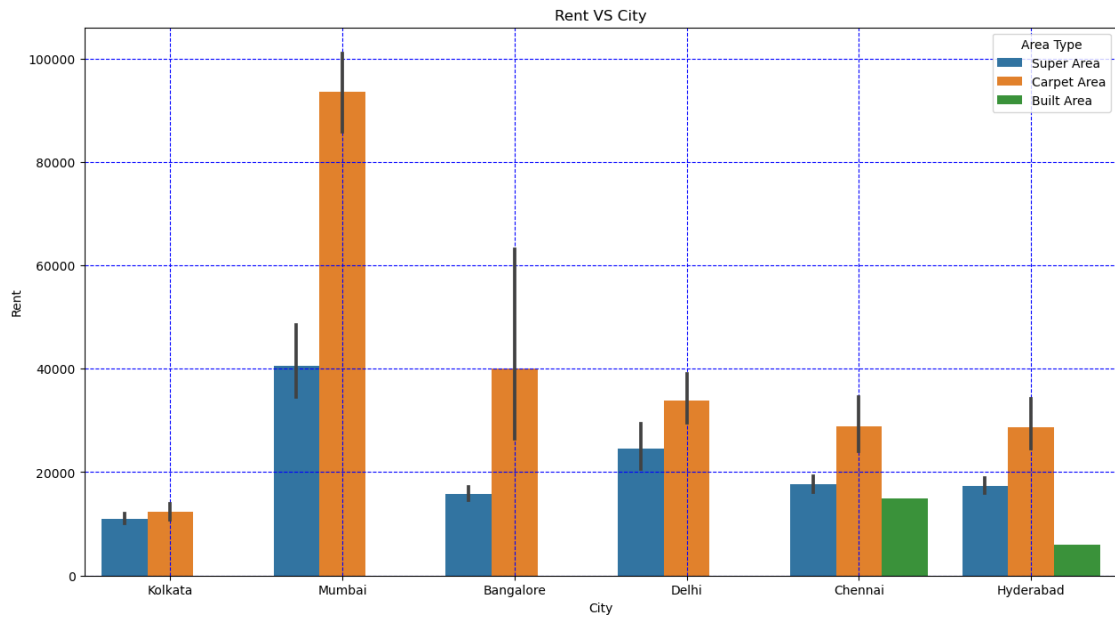


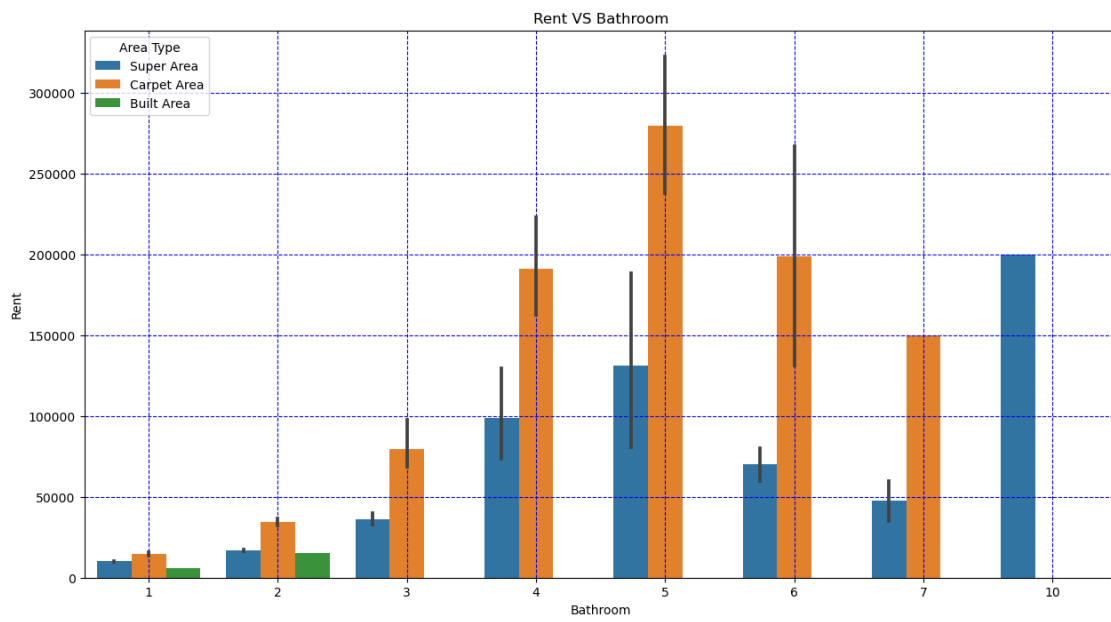
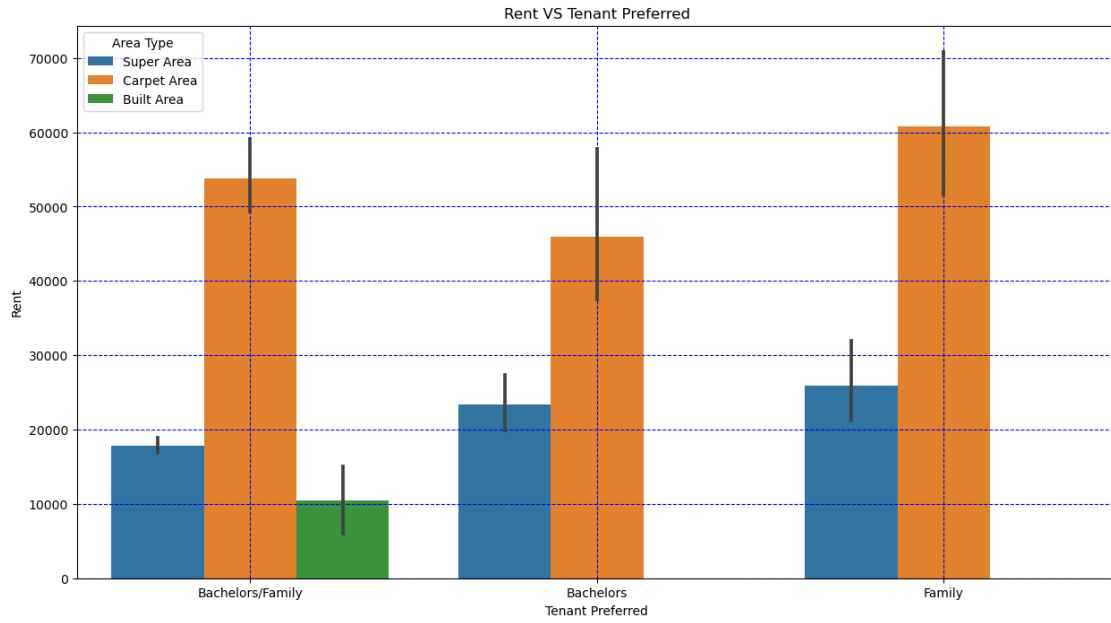
```
[24]: target_column = 'Rent'
sub_columns = [z for z in df.columns if z != target_column]
for sub_column in sub_columns:
    plt.figure(figsize=(15,8))
    sns.barplot(x = sub_column, y=target_column, data=df, hue = 'Area Type')
    plt.title(f'{target_column} VS {sub_column}')
    plt.grid(True, linestyle='--', color='blue')
    plt.show()
```

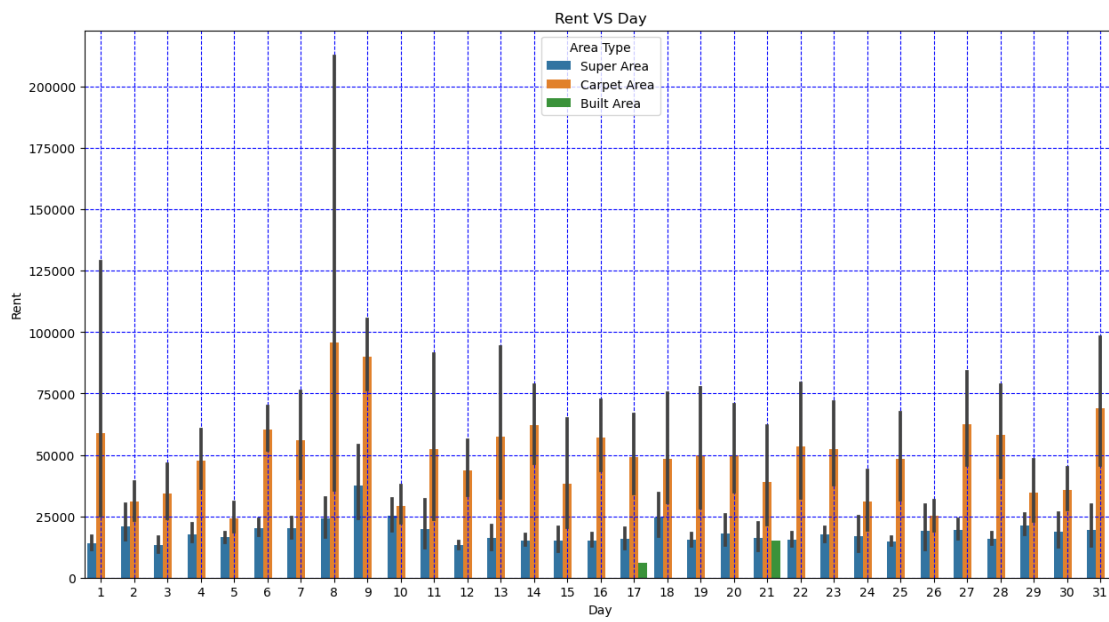
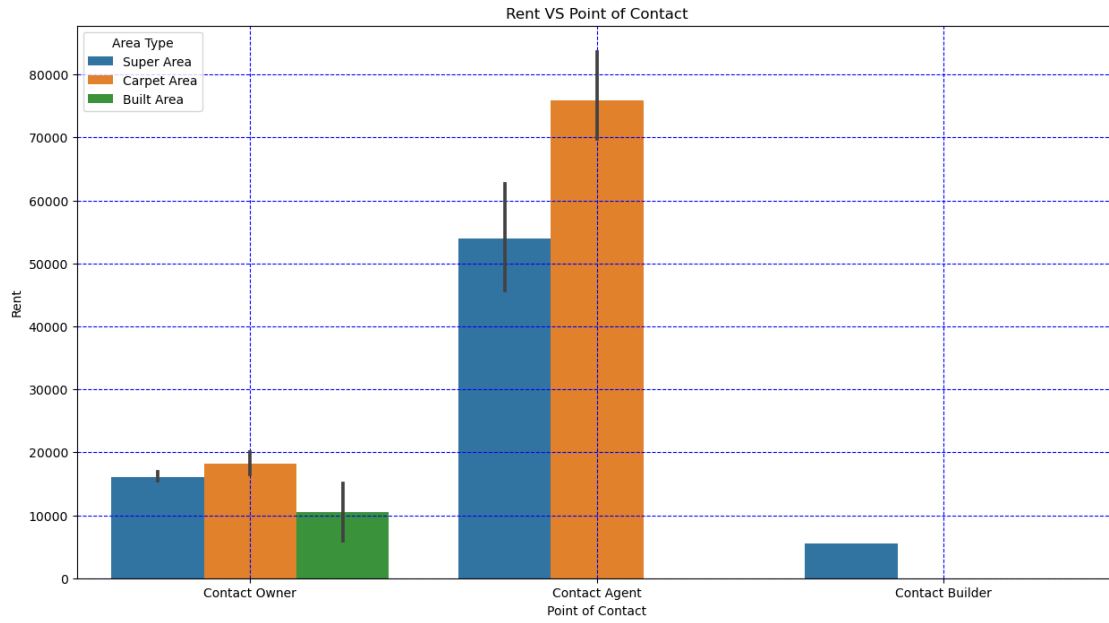


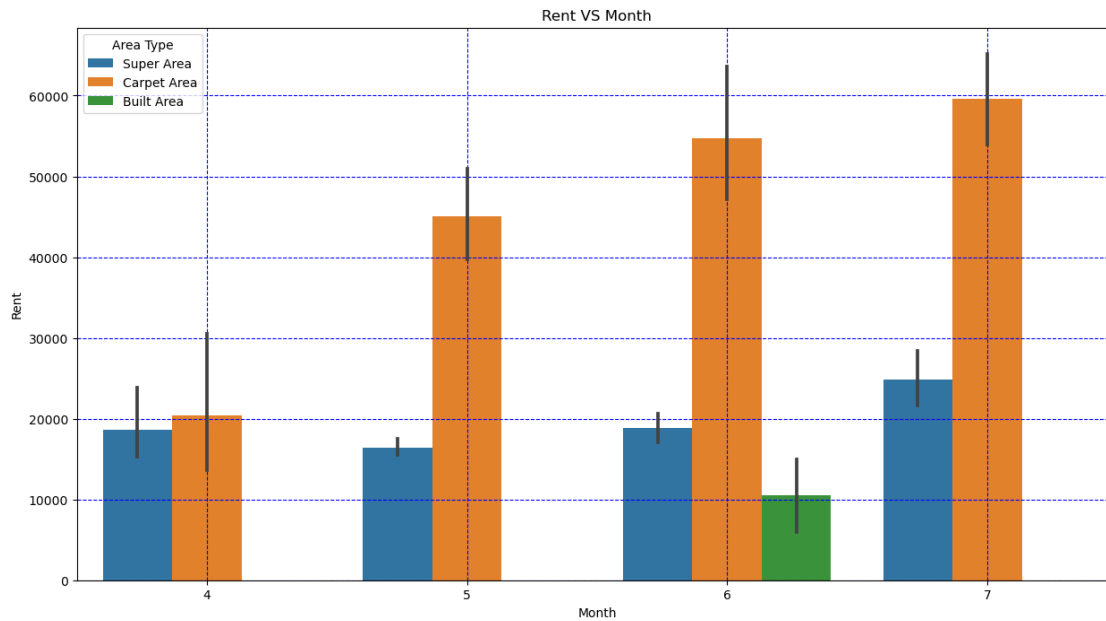










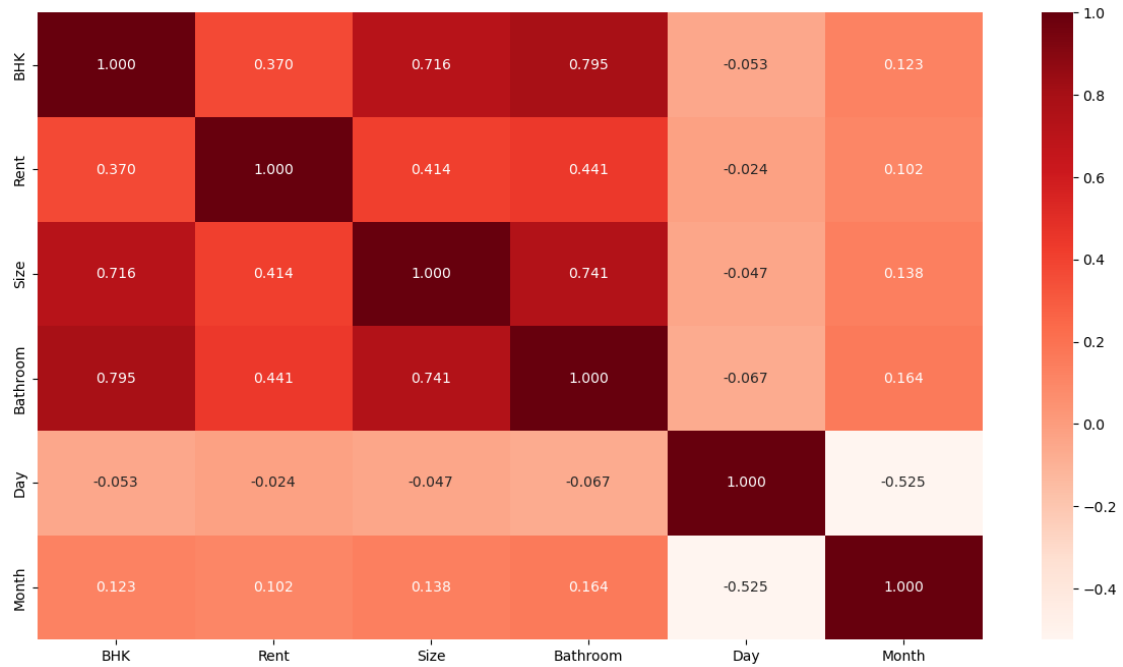


```
[28]: plt.figure(figsize=(15,8))
sns.heatmap(df.corr(), annot=True, cmap='Reds', fmt='.3f')
plt.show()
```

C:\Users\Admin\AppData\Local\Temp\ipykernel\_7396\2229025032.py:2: FutureWarning:  
The default value of numeric\_only in DataFrame.corr is deprecated. In a future  
version, it will default to False. Select only valid columns or specify the  
value of numeric\_only to silence this warning.

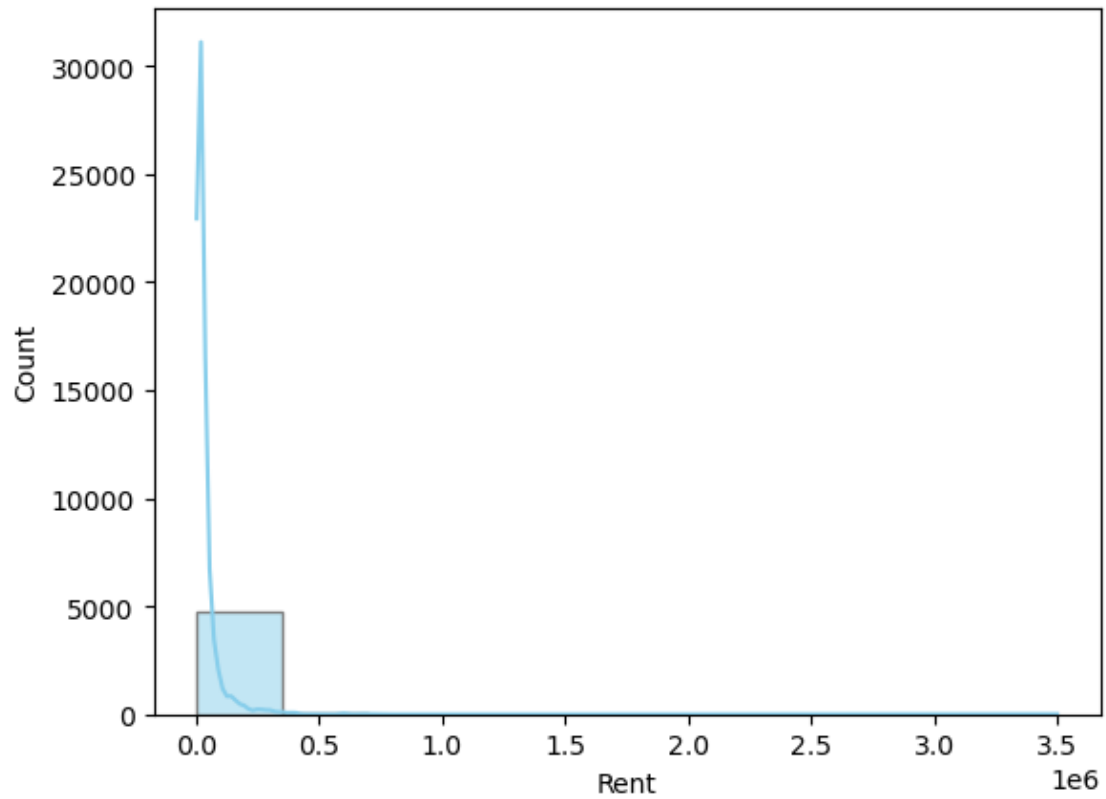
```
sns.heatmap(df.corr(), annot=True, cmap='Reds', fmt='.3f')
```



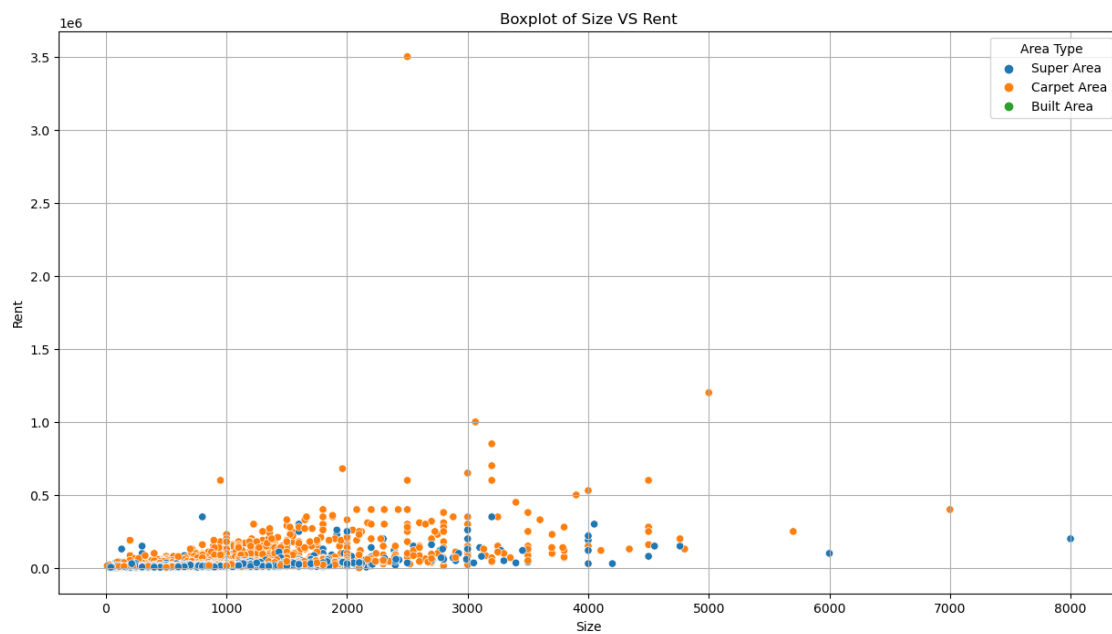
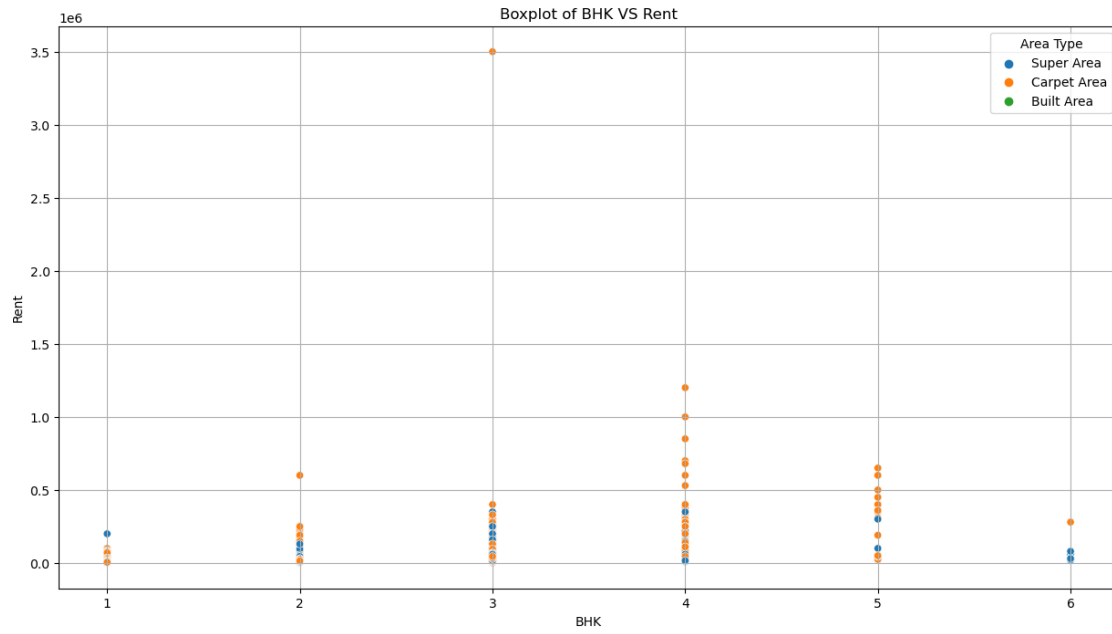


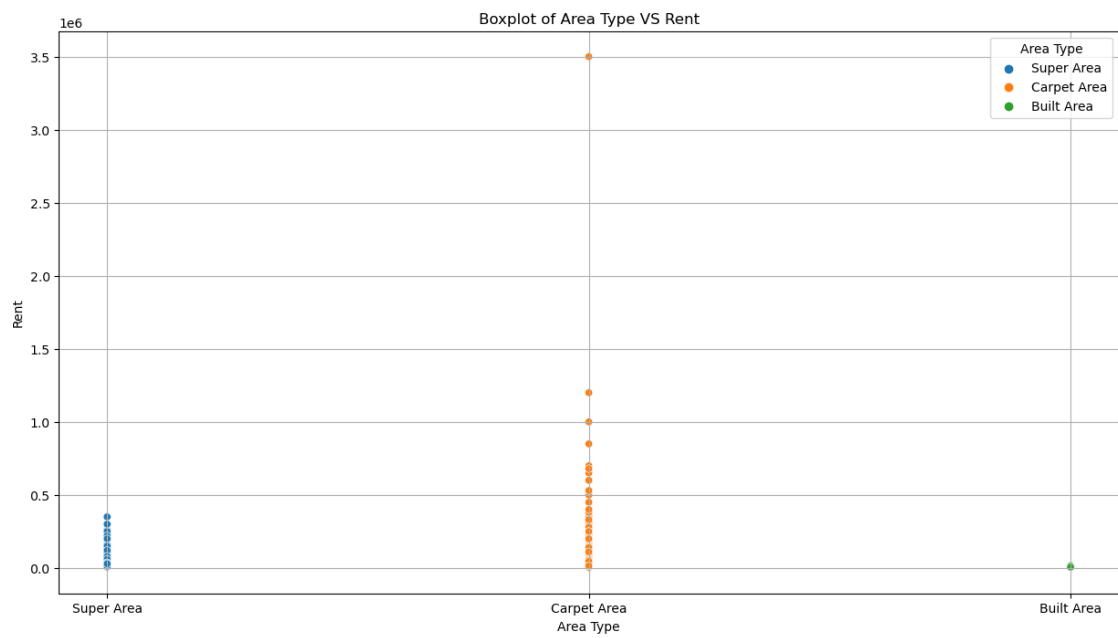
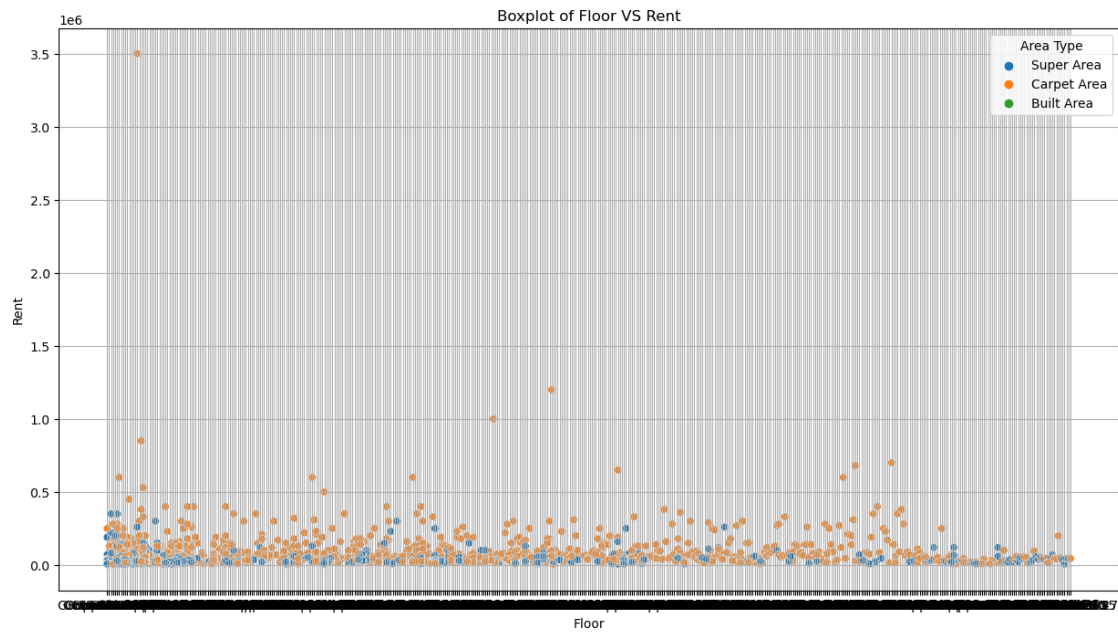
```
[30]: sns.histplot(df['Rent'], kde=True, bins=10, color='skyblue', edgecolor='grey')
```

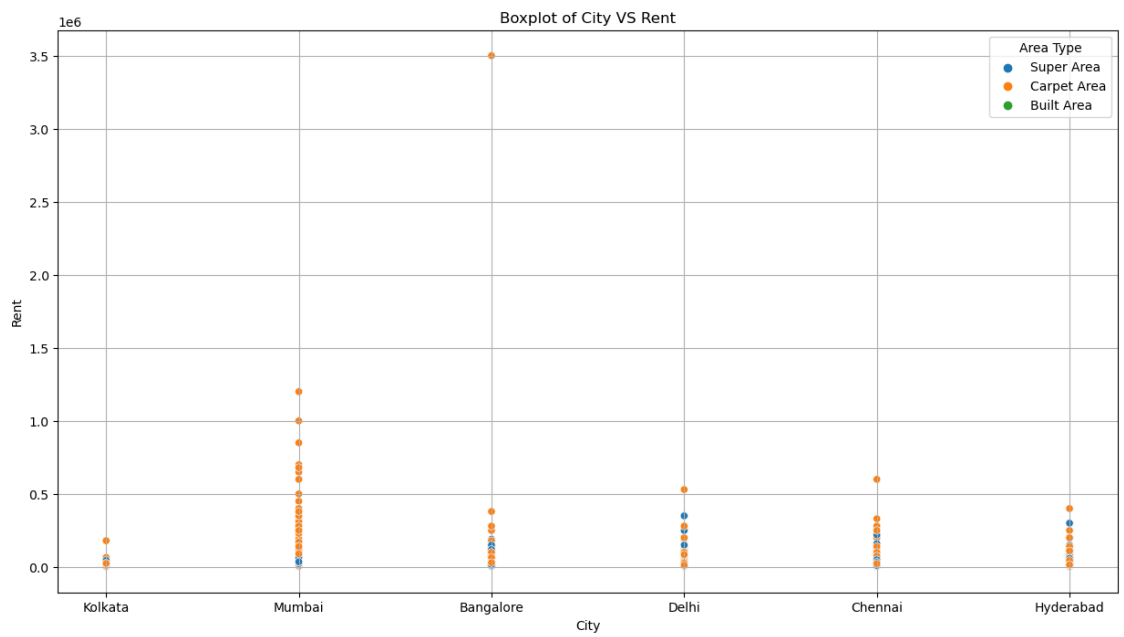
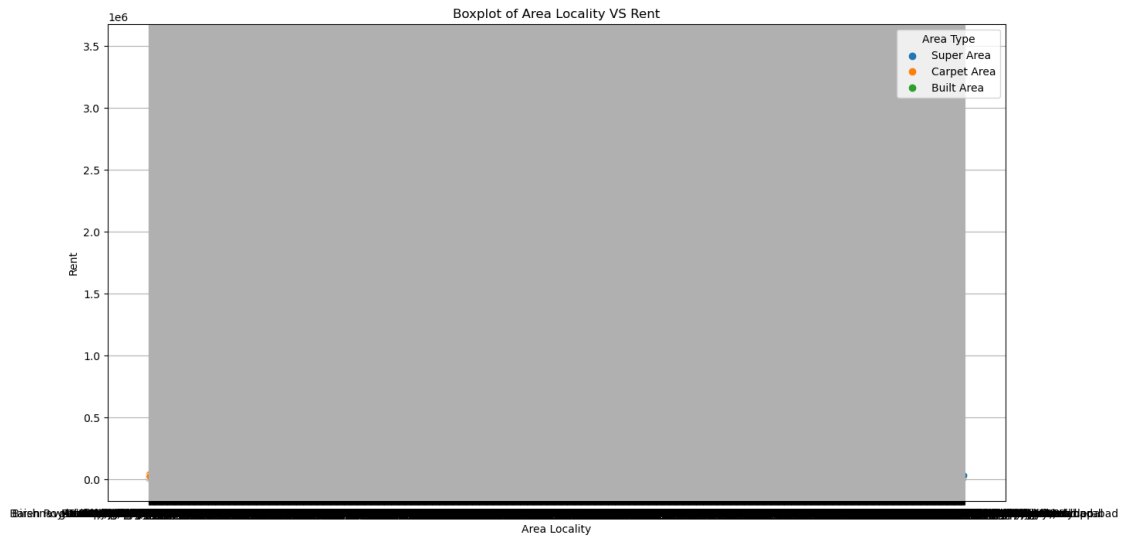
```
[30]: <Axes: xlabel='Rent', ylabel='Count'>
```

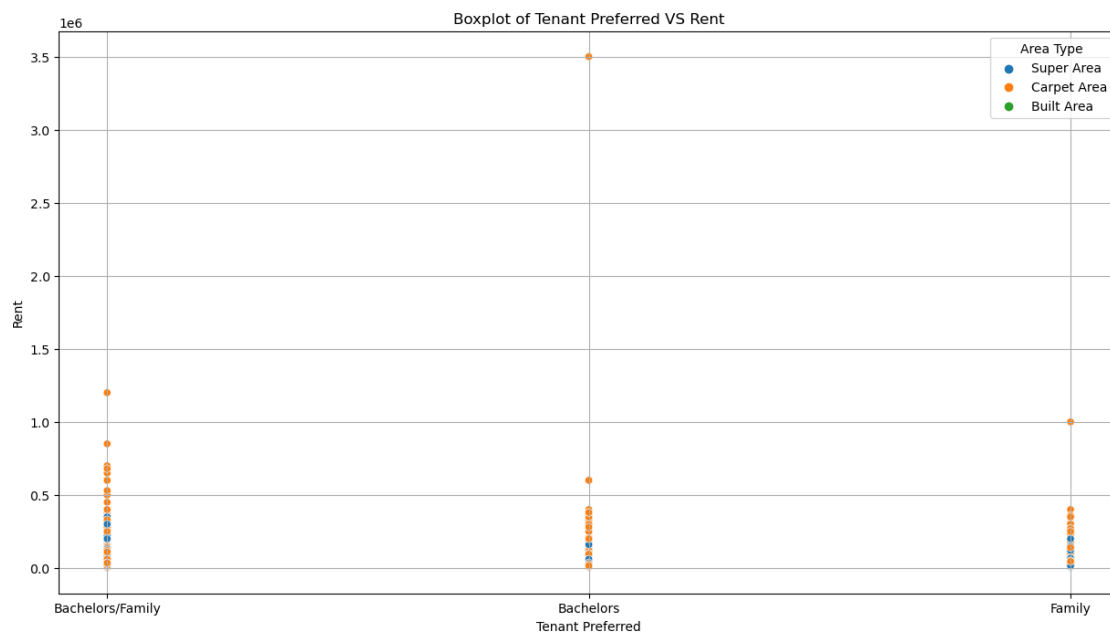
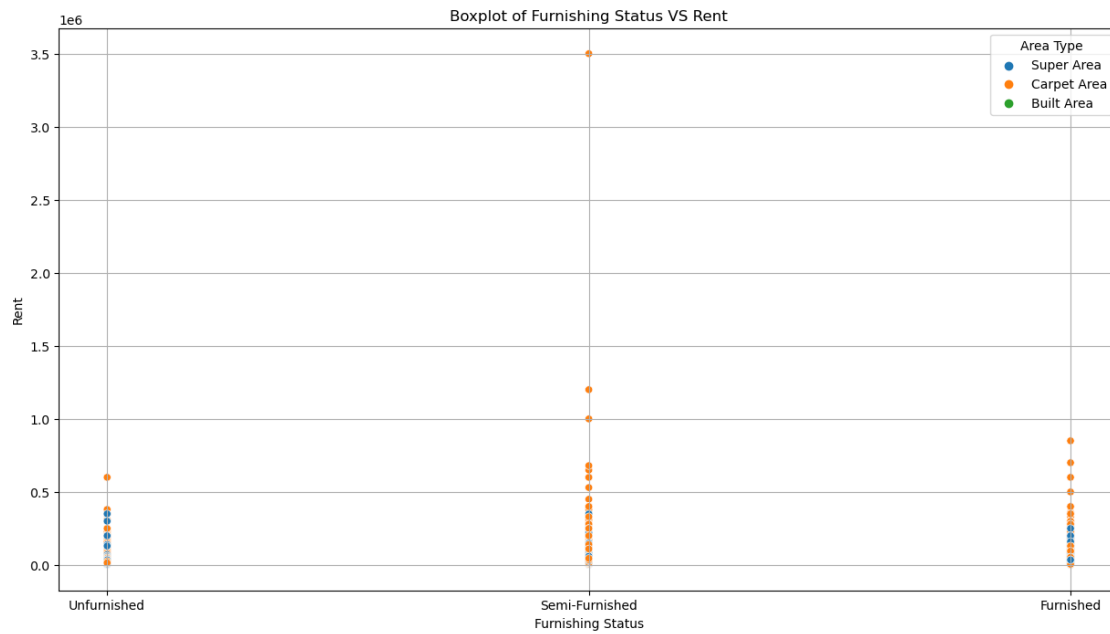


```
[42]: target_column = 'Rent'
sub_columns = [m for m in df.columns if m != target_column]
for sub_column in sub_columns:
    plt.figure(figsize = (15,8))
    sns.scatterplot(x=sub_column, y = 'Rent', data=df, hue='Area Type')
    plt.title(f'Boxplot of {sub_column} VS {target_column}')
    plt.grid()
    plt.show()
```

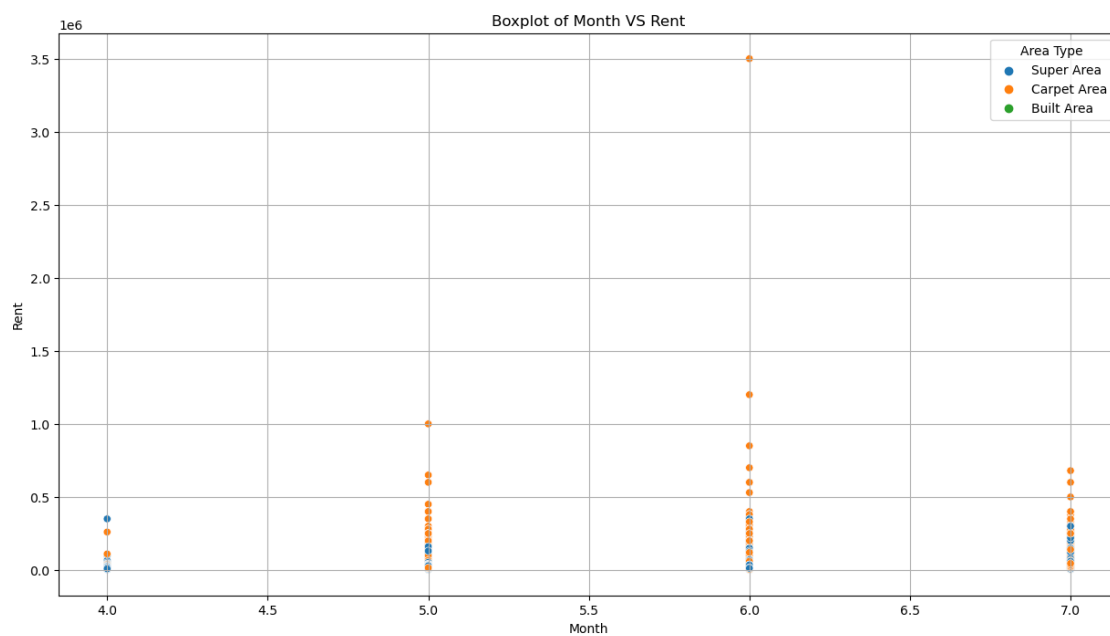
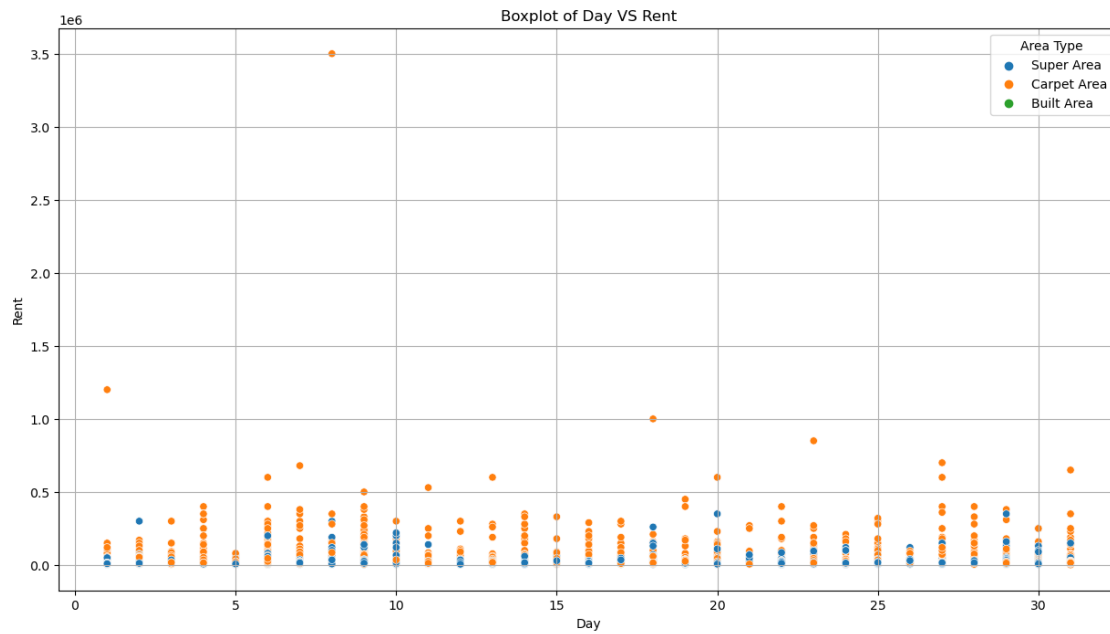








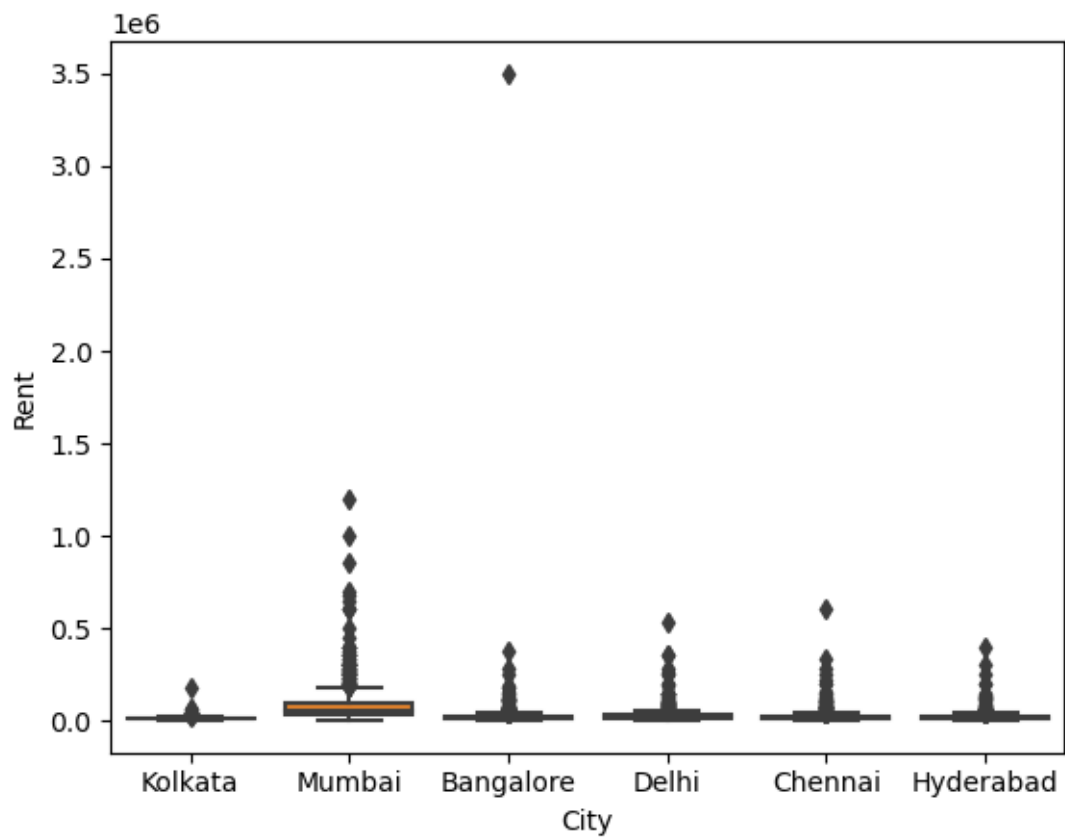




```
[37]: sns.boxplot(x='City', y = 'Rent', data=df)
```

```
[37]: <Axes: xlabel='City', ylabel='Rent'>
```





```
[40]: sns.pairplot(df)
```

```
[40]: <seaborn.axisgrid.PairGrid at 0x1a06cd61490>
```



```
[43]: df.head()
```

```
[43]:
```

	BHK	Rent	Size	Floor	Area Type	Area Locality \
0	2	10000	1100	Ground out of 2	Super Area	Bandel
1	2	20000	800	1 out of 3	Super Area	Phool Bagan, Kankurgachi
2	2	17000	1000	1 out of 3	Super Area	Salt Lake City Sector 2
3	2	10000	800	1 out of 2	Super Area	Dum Dum Park
4	2	7500	850	1 out of 2	Carpet Area	South Dum Dum

	City	Furnishing Status	Tenant Preferred	Bathroom	Point of Contact \
0	Kolkata	Unfurnished	Bachelors/Family	2	Contact Owner
1	Kolkata	Semi-Furnished	Bachelors/Family	1	Contact Owner
2	Kolkata	Semi-Furnished	Bachelors/Family	1	Contact Owner

3	Kolkata	Unfurnished	Bachelors/Family	1	Contact Owner
4	Kolkata	Unfurnished	Bachelors	1	Contact Owner

	Day	Month
0	18	5
1	13	5
2	16	5
3	4	7
4	9	5

```
[65]: column_plot = ['BHK', 'Rent', 'Size', 'Floor', 'Area Type', 'Area Locality',
                    'City', 'Furnishing Status', 'Tenant Preferred', 'Bathroom', 'Point_
                    of Contact', 'Day', 'Month']
column_plot
```

```
[65]: ['BHK',
        'Rent',
        'Size',
        'Floor',
        'Area Type',
        'Area Locality',
        'City',
        'Furnishing Status',
        'Tenant Preferred',
        'Bathroom',
        'Point of Contact',
        'Day',
        'Month']
```

```
[76]: #create number of columns for subplots
num_columns=2
num_rows = (len(column_plot)//num_columns) + (len(column_plot) % num_columns >
0)

#create subplots
fig, axes = plt.subplots(num_rows, num_columns, figsize = (15,5 * num_rows))
fig.suptitle('Pie Charts of column Distributions',y=1.02)

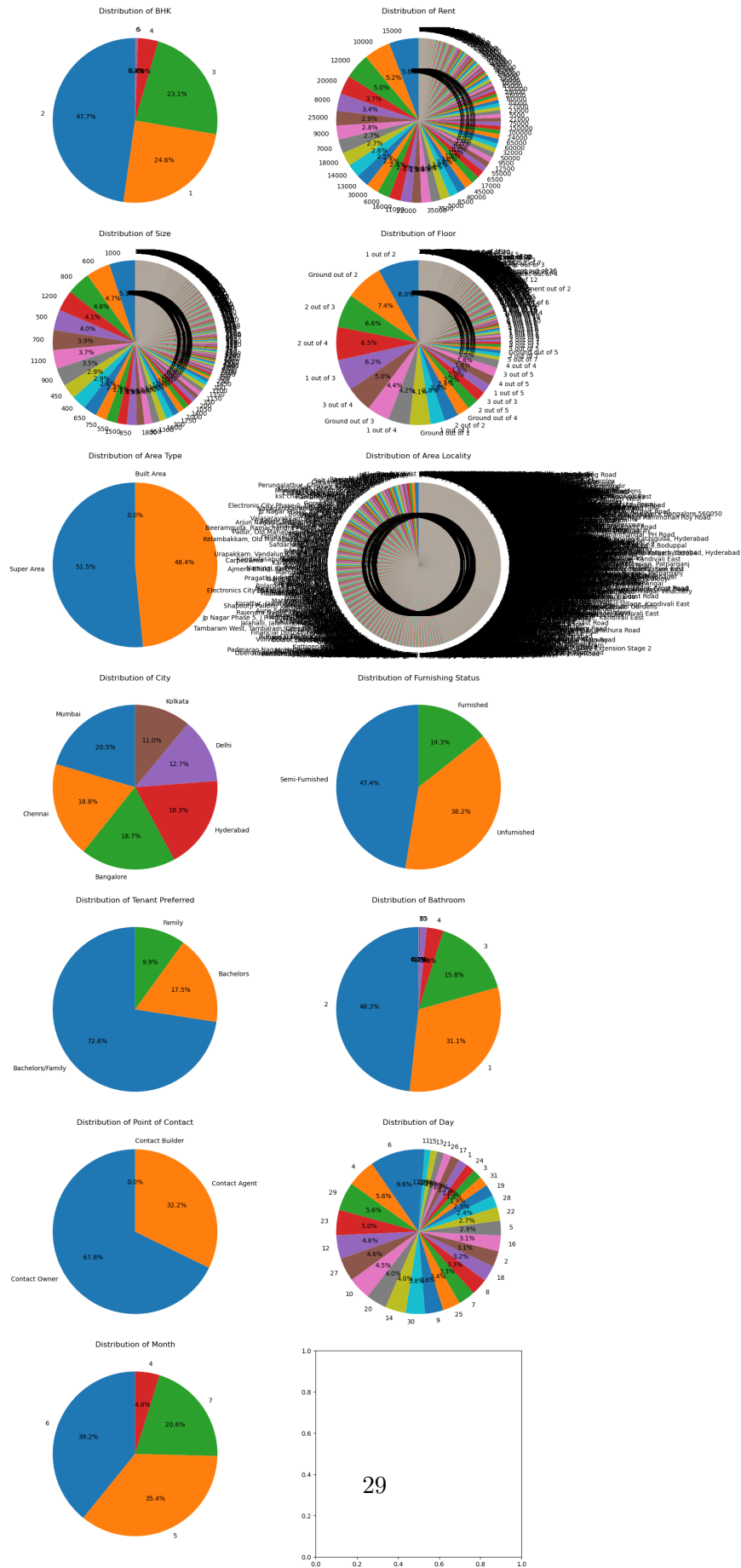
# Easier for indexing
axes = axes.flatten()

for x, column in enumerate(column_plot):
    column_count = df[column].value_counts()

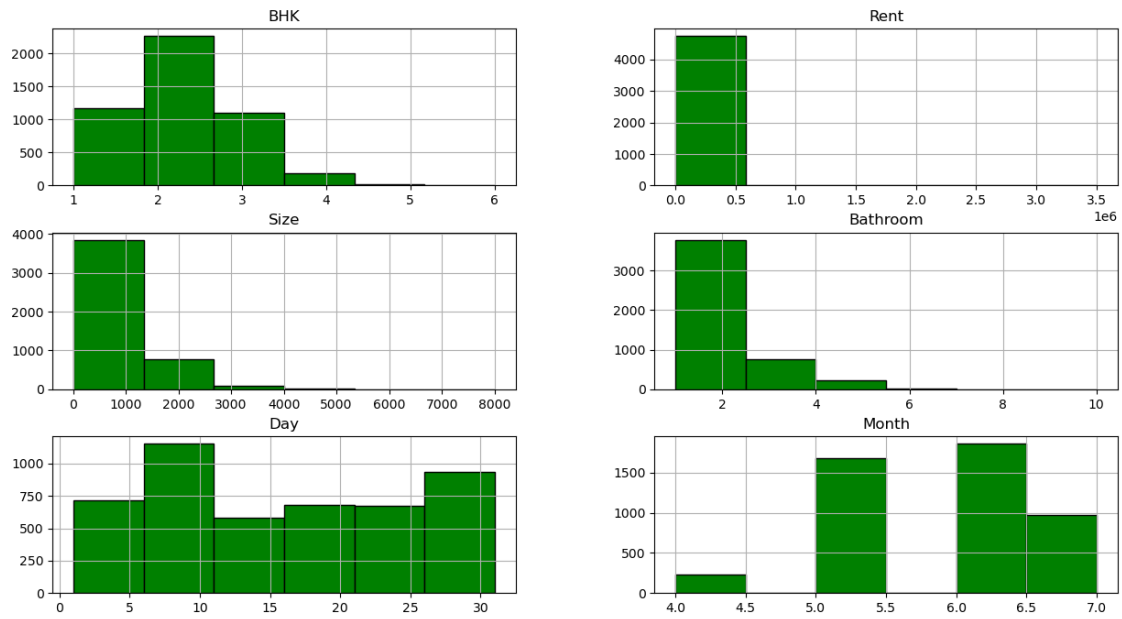
    sns.color_palette('pastel')
```

```
    axes[x].pie(column_count, labels = column_count.index, autopct = '%1.1f%%',  
↪startangle=90)  
    axes[x].set_title(f'Distribution of {column}')  
  
plt.tight_layout()  
plt.show()
```

Pie Charts of column Distributions



```
[79]: df.hist(bins=6,figsize=(15,8), color='green', edgecolor='black')  
plt.show()
```



```
[ ]:
```