

# Dubai Rental Dataset Analysis

## Purpose of work:

- To investigate trends, patterns and insights from the dataset on rentals in the city of Dubai and
- To predict rental values based on features in the dataset using machine learning models.

[Dataset \(<https://www.kaggle.com/datasets/azharsaleem/real-estate-goldmine-dubai-uae-rental-market/data>\)](https://www.kaggle.com/datasets/azharsaleem/real-estate-goldmine-dubai-uae-rental-market/data)

## Dataset Overview

Each entry in the dataset represents a rental property listing with details about the property's features, rental terms, and location specifics. This primary and unique dataset is designed for analysis and can be used to generate insights into the rental market dynamics of the UAE.

## Columns Description

Address: Full address of the property.

Rent: The annual rent price in AED.

Beds: Number of bedrooms in the property.

Baths: Number of bathrooms in the property.

Type: Type of property (e.g., Apartment, Villa, Penthouse).

Area\_in\_sqft: Total area of the property in square feet.

Rent\_per\_sqft: Rent price per square foot, calculated as Rent divided by Area\_in\_sqft.

Rent\_category: Categorization of the rent price (Low, Medium, High) based on thresholds.

Frequency: Rental payment frequency, which is consistently 'Yearly'.

Furnishing: Furnishing status of the property (Furnished, Unfurnished).

Purpose: The purpose of the listing, typically 'For Rent'.

Posted\_date: The date the property was listed for rent.

Age\_of\_listing\_in\_days: The number of days the listing has been active since it was posted.

Location: A more specific location within the city where the property is located.

City: City in which the property is situated.

Latitude, Longitude: Geographic coordinates of the property.

## 1.0 Import Libraries and Data

## 2.0 Data Integrity Checks

## 3.0 Exploratory Data Analysis

### 3.1 Property type, rental category, furnishing, bedrooms and bathrooms

### 3.2 Rents over time

### 3.3 Average rents over time

### 3.4 Age of listing

### 3.5 Location

### 3.6 Area

### 3.7 Relationships

## 4.0 Prediction

### 4.1 Data preparation

### 4.2 Linear Regression

### 4.3 Other Models

## 1.0 Import Libraries and Data

```
In [996]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 from scipy import stats
6 import folium
7 import warnings
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.preprocessing import LabelEncoder
10 from sklearn.model_selection import train_test_split
11 from sklearn.linear_model import LinearRegression
12 from sklearn import neighbors
13 from sklearn.ensemble import RandomForestClassifier
14 from sklearn.metrics import r2_score
15 warnings.filterwarnings("ignore")
```

```
In [287]: 1 df = pd.read_csv(r'C:\Users\imoge\AllMLProjects\Data\dubai_properties.csv')
```

```
In [288]: 1 print(df.shape)
```

```
(73742, 17)
```

```
In [289]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73742 entries, 0 to 73741
Data columns (total 17 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Address          73742 non-null   object 
 1   Rent              73742 non-null   int64  
 2   Beds              73742 non-null   int64  
 3   Baths             73742 non-null   int64  
 4   Type              73742 non-null   object 
 5   Area_in_sqft     73742 non-null   int64  
 6   Rent_per_sqft    73742 non-null   float64
 7   Rent_category    73742 non-null   object 
 8   Frequency         73742 non-null   object 
 9   Furnishing        73742 non-null   object 
 10  Purpose            73742 non-null   object 
 11  Posted_date       73742 non-null   object 
 12  Age_of_listing_in_days 73742 non-null   int64  
 13  Location           73742 non-null   object 
 14  City               73742 non-null   object 
 15  Latitude           73023 non-null   float64
 16  Longitude          73023 non-null   float64
dtypes: float64(3), int64(5), object(9)
memory usage: 9.6+ MB
```

## 2.0 Data Integrity Checks

```
In [290]: 1 # Check the value counts by city
2 df['City'].value_counts()
```

```
Out[290]: Dubai           34250
Abu Dhabi        23324
Sharjah          9516
Ajman            4704
Al Ain           1040
Ras Al Khaimah   816
Umm Al Quwain    65
Fujairah         27
Name: City, dtype: int64
```

```
In [291]: 1 # Split out just the Dubai results
2 dubai = df[df['City']=='Dubai']
3 print(dubai.shape)
```

(34250, 17)

```
In [292]: 1 dubai.isnull().sum()
```

```
Out[292]: Address      0
Rent        0
Beds        0
Baths       0
Type        0
Area_in_sqft  0
Rent_per_sqft 0
Rent_category 0
Frequency    0
Furnishing   0
Purpose      0
Posted_date  0
Age_of_listing_in_days 0
Location     0
City         0
Latitude     31
Longitude    31
dtype: int64
```

```
In [293]: 1 # We can see all the missing items relate to one location
2 dubai[dubai['Latitude'].isnull()]['Location'].value_counts()
```

```
Out[293]: Dubai Science Park    31
Name: Location, dtype: int64
```

```
In [294]: 1 # Lets replace these values with the Latitude and longitude for that area
2 latitude = 25.0745
3 longitude = 55.2396
4
5 dubai['Latitude'].fillna(latitude, inplace = True)
6 dubai['Longitude'].fillna(longitude, inplace = True)
```

```
In [295]: 1 # Do we have any duplicates?
2 dubai.duplicated().value_counts()
```

```
Out[295]: False    34250
dtype: int64
```

We don't need the purpose, frequency and city columns have only a single value so are not useful

```
In [296]: 1 # Drop columns we dont need
2 dubai.drop(columns = ['Purpose', 'Frequency', 'City'], axis = 1, inplace = True)
3 dubai.head()
```

Out[296]:

|       | Address  | Rent   | Beds | Baths | Type      | Area_in_sqft | Rent_per_sqft | Rent_category | Furnishing  | Posted_date | Age_of_listing |
|-------|--|--------|------|-------|-----------|--------------|---------------|---------------|-------------|-------------|----------------|
| 29068 | Binghatti Heights,<br>JVC District<br>10,<br>Jumeirah V... | 125000 | 2    | 2     | Apartment | 1145         | 109.170306    | Medium        | Unfurnished | 2024-02-15  |                |
| 29069 | Seasons Community,<br>JVC District<br>15,<br>Jumeirah V... | 50000  | 1    | 2     | Apartment | 655          | 76.335878     | Low           | Unfurnished | 2024-03-21  |                |
| 29070 | Autumn 1 Block B,<br>Autumn Seasons Community,<br>J...     | 90000  | 1    | 2     | Apartment | 896          | 100.446429    | Medium        | Furnished   | 2024-04-08  |                |
| 29071 | Socio Tower A,<br>Socio,<br>Dubai Hills Estate,<br>Dubai   | 125000 | 2    | 1     | Apartment | 720          | 173.611111    | Medium        | Unfurnished | 2024-03-11  |                |
| 29072 | Eleganz by Danube,<br>JVC District<br>14,<br>Jumeirah V... | 105000 | 1    | 2     | Apartment | 965          | 108.808290    | Medium        | Furnished   | 2024-02-23  |                |



## Rent

```
In [298]: 1 # Descriptive Statistics
2 dubai.describe()
```

Out[298]:

|              | Rent         | Beds         | Baths        | Area_in_sqft  | Rent_per_sqft | Age_of_listing_in_days | Latitude     | Longitude    |
|--------------|--------------|--------------|--------------|---------------|---------------|------------------------|--------------|--------------|
| <b>count</b> | 3.425000e+04 | 34250.000000 | 34250.000000 | 34250.000000  | 34250.000000  | 34250.000000           | 34250.000000 | 34250.000000 |
| <b>mean</b>  | 2.133664e+05 | 1.971212     | 2.081518     | 1831.808321   | 132.253717    | 63.912263              | 25.115431    | 55.243365    |
| <b>std</b>   | 4.272489e+05 | 1.395483     | 0.561300     | 3119.024048   | 70.329882     | 55.176208              | 0.072873     | 0.079815     |
| <b>min</b>   | 0.000000e+00 | 0.000000     | 1.000000     | 74.000000     | 0.000000      | 12.000000              | 24.839901    | 54.998449    |
| <b>25%</b>   | 8.500000e+04 | 1.000000     | 2.000000     | 754.000000    | 86.614173     | 27.000000              | 25.058657    | 55.170057    |
| <b>50%</b>   | 1.450000e+05 | 2.000000     | 2.000000     | 1163.000000   | 120.000000    | 46.000000              | 25.098721    | 55.263874    |
| <b>75%</b>   | 2.300000e+05 | 3.000000     | 2.000000     | 1930.750000   | 159.997167    | 81.000000              | 25.186684    | 55.288876    |
| <b>max</b>   | 5.500000e+07 | 12.000000    | 11.000000    | 210254.000000 | 2182.044888   | 1131.000000            | 25.300533    | 55.569534    |

The most expensive rental is 55 million dirham per year (approx \$15m). The least expensive is 0. This is a bit odd as we would expect a value for all properties.

There is right skew in the rental data with the mean higher than the median and there are likely outliers (we will look at this again in a bit)

```
In [299]: 1 # Places with zero rent
2 dubai[dubai['Rent']==0]
```

Out[299]:

|  |       | Address   | Rent | Beds | Baths | Type      | Area_in_sqft | Rent_per_sqft | Rent_category | Furnishing  | Posted_date | Age_of_listir |
|--|-------|---|------|------|-------|-----------|--------------|---------------|---------------|-------------|-------------|---------------|
|  | 32240 | Al Barsha South 2, Al Barsha South, Al Barsha,... | 0    | 5    | 2     | Villa     | 10000        | 0.0           | Low           | Unfurnished | 2024-03-26  |               |
|  | 32241 | Port Saeed, Deira, Dubai                          | 0    | 1    | 2     | Apartment | 750          | 0.0           | Low           | Unfurnished | 2024-03-06  |               |
|  | 32242 | Port Saeed, Deira, Dubai                          | 0    | 1    | 2     | Apartment | 850          | 0.0           | Low           | Unfurnished | 2024-03-06  |               |
|  | 32243 | Park Gate Residence A, Park Gate Residence, Al... | 0    | 2    | 2     | Apartment | 1554         | 0.0           | Low           | Furnished   | 2024-01-30  |               |
|  | 32631 | Al Barsha South 2, Al Barsha South, Al Barsha,... | 0    | 5    | 2     | Villa     | 12500        | 0.0           | Low           | Unfurnished | 2024-03-04  |               |
|  | 32632 | Abu Hail, Deira, Dubai                            | 0    | 5    | 2     | Villa     | 4500         | 0.0           | Low           | Unfurnished | 2023-12-13  |               |
|  | 32633 | Al Barsha 3, Al Barsha, Dubai                     | 0    | 1    | 2     | Apartment | 750          | 0.0           | Low           | Furnished   | 2024-03-21  |               |
|  | 34160 | Al Khudrawi, Shoreline Apartments, Palm Jumeir... | 0    | 2    | 2     | Apartment | 1551         | 0.0           | Low           | Unfurnished | 2024-03-05  |               |
|  | 35389 | Serenia Residences West Wing, Serenia Residenc... | 0    | 4    | 2     | Apartment | 2885         | 0.0           | Low           | Furnished   | 2024-03-18  |               |
|  | 56079 | 1 Residences, Wasl 1, Al Kifaf, Bur Dubai, Dubai  | 0    | 2    | 2     | Apartment | 1439         | 0.0           | Low           | Unfurnished | 2024-03-10  |               |
|  | 56080 | Abu Keibal, Shoreline Apartments, Palm Jumeira... | 0    | 3    | 2     | Apartment | 2138         | 0.0           | Low           | Unfurnished | 2024-03-25  |               |
|  | 56081 | Garden Homes Frond K, Garden Homes Palm Jumeir... | 0    | 4    | 2     | Villa     | 5000         | 0.0           | Low           | Unfurnished | 2024-03-20  |               |
|  | 56082 | Marina Residences 5, Marina Residences, Palm J... | 0    | 2    | 2     | Apartment | 3996         | 0.0           | Low           | Furnished   | 2023-11-17  |               |
|  | 56083 | Al Jafiliya, Dubai                                | 0    | 5    | 2     | Villa     | 7000         | 0.0           | Low           | Unfurnished | 2023-11-24  |               |

It is not clear why these properties are listed at zero rent. Perhaps there is some kind of maintenance charge or other charge and/or they are less desirable. There are also some with a peppercorn rent of just 1 AED. We will drop these from our data as we don't know why this is the case and it is not likely to add value to the analysis to include them.

```
In [300]: 1 # Drop zero rentals and a few posted at 1 AED  
2 dubai = dubai[dubai['Rent']>1]  
3 dubai.shape
```

Out[300]: (34236, 14)

We should probably remove the high rent outliers also

```
In [516]: 1 # Create a zscore column and exclude those above and below 3 standard deviations from the mean  
2 dubai['RentZScore'] = stats.zscore(dubai['Rent'])
```

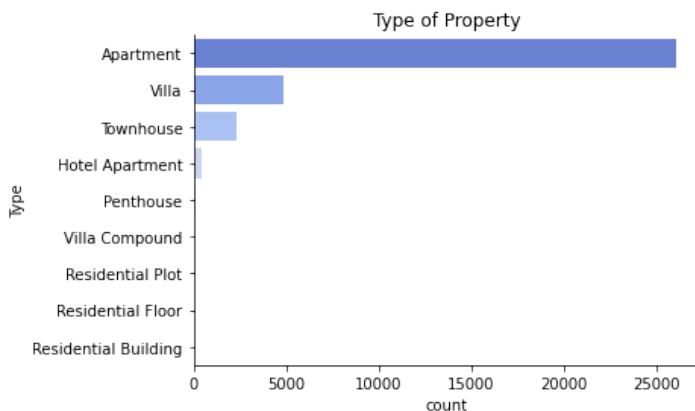
```
In [526]: 1 # Exclude outliers  
2 dubai = dubai[(dubai['RentZScore']<=3) & (dubai['RentZScore']>-3)]
```

## 3.0 Exploratory Data Analysis

We have a mix of categorical and numerical variables that we can investigate both separately and in terms of relationships between them.

### 3.1 Property type, rental category, furnishing, bedrooms and bathrooms

```
In [527]: 1 # Plot the type of property distribution  
2 ax = sns.countplot(data = dubai, y = 'Type', palette="coolwarm" )  
3 plt.title("Type of Property")  
4 ax.spines.right.set_visible(False)  
5 ax.spines.top.set_visible(False);
```



```
In [528]: 1 dubai['Type'].value_counts(normalize = True)*100
```

Out[528]:

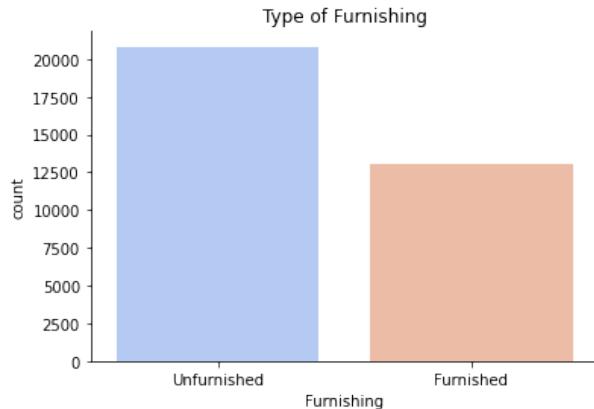
|                      |           |
|----------------------|-----------|
| Apartment            | 76.993657 |
| Villa                | 14.238088 |
| Townhouse            | 6.944977  |
| Hotel Apartment      | 1.327629  |
| Penthouse            | 0.418941  |
| Villa Compound       | 0.050155  |
| Residential Building | 0.011801  |
| Residential Floor    | 0.008851  |
| Residential Plot     | 0.005901  |

Name: Type, dtype: float64

Some 76% of properties are apartments with a further 15% being villas

In [529]:

```
1 # Plot the furnishing distribution
2 ax = sns.countplot(data = dubai, x = 'Furnishing', palette="coolwarm" )
3 plt.title("Type of Furnishing")
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



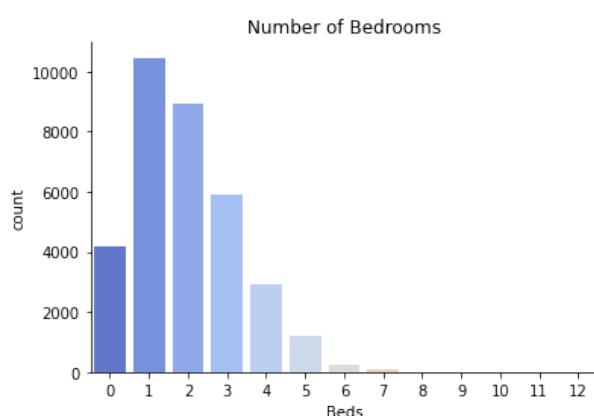
In [530]:

```
1 # Plot the rent category distribution
2 ax = sns.countplot(data = dubai, x = 'Rent_category', palette="coolwarm" )
3 plt.title("Rent Category")
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



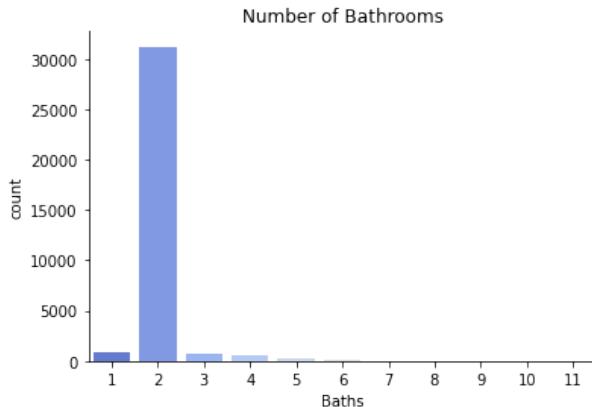
In [531]:

```
1 # Plot the number of bedrooms distribution
2 ax = sns.countplot(data = dubai, x = 'Beds', palette="coolwarm" )
3 plt.title("Number of Bedrooms")
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



In [532]:

```
1 # Plot the number of bathrooms distribution
2 ax = sns.countplot(data = dubai, x = 'Baths', palette="coolwarm" )
3 plt.title("Number of Bathrooms")
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



## Summary

Most properties are high rental unfurnished apartments with one bedroom and two bathrooms.

## 3.2 Rentals Over Time

In [533]:

```
1 # Convert the posted date to a datetime object
2 dubai['Posted_date'] = pd.to_datetime(dubai['Posted_date'])
3
4 # Split out the year and month for analysis
5 dubai['Year'] = dubai['Posted_date'].dt.year
6 dubai['Month'] = dubai['Posted_date'].dt.month
7 dubai['MonthName'] = dubai['Posted_date'].dt.month_name()
```

In [534]:

```
1 # Get the minimum and maximum posting date
2 print(dubai['Posted_date'].min())
3 print(dubai['Posted_date'].max())
```

2021-03-17 00:00:00  
2024-04-09 00:00:00

In [535]:

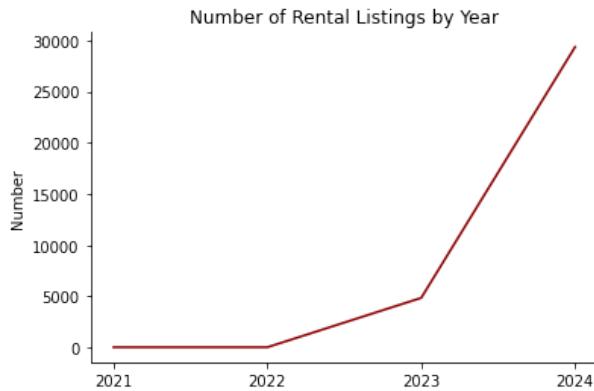
```
1 # Rentals by year
2 years2 = dubai.groupby('Year',as_index = False)[['Rent']].count()
3 years2.columns = ['Year', 'Number']
4 years2
```

Out[535]:

|   | Year | Number |
|---|------|--------|
| 0 | 2021 | 15     |
| 1 | 2022 | 15     |
| 2 | 2023 | 4766   |
| 3 | 2024 | 29099  |

In [536]:

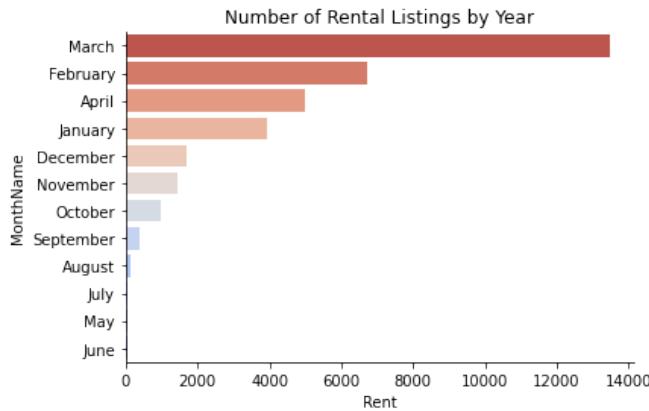
```
1 # Create a Line plot
2 xvals = ['2021','2022','2023','2024']
3 ax = sns.lineplot(data=years, x=xvals, y="Number", color = 'darkred')
4 plt.title("Number of Rental Listings by Year")
5 ax.spines.right.set_visible(False)
6 ax.spines.top.set_visible(False);
```



The data covers just over 3 years from mid March 2021 to just into April 2024. We can see the number of rentals listed has increased significantly between 2023 and 2024 even though only 3 months of the latest year are included. Given this extremely large increase, it might be worth checking the original dataset to see how it was collected and to ensure that we have full information for the prior years and checking the figures for the rest of 2024 when they become available.

In [537]:

```
1 # Number of rentals by Month
2 years3 = dubai.groupby('MonthName',as_index = False)[['Rent']].count().sort_values(by = 'Rent',ascending
3
4 # Create a line plot
5 ax = sns.barplot(data=years3, y='MonthName', x="Rent",palette = 'coolwarm_r')
6 plt.title("Number of Rental Listings by Year")
7 ax.spines.right.set_visible(False)
8 ax.spines.top.set_visible(False);
```



More properties are listed in March and February and very few in the summer period of May through to August

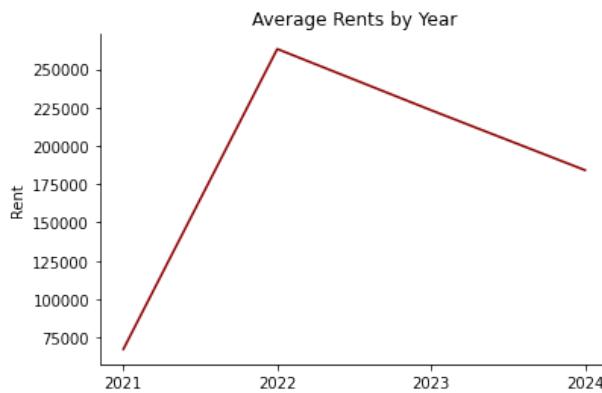
### 3.3 Average Rents over Time

In [538]:

```
1 print(dubai['Rent'].mean())
2 print(dubai['Rent'].median())
```

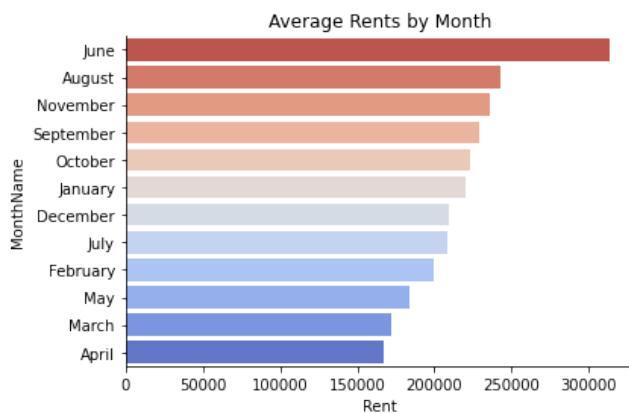
```
189587.1648030683
142000.0
```

```
In [539]: 1 # Rents by year
2 years2 = dubai.groupby('Year',as_index = False)[['Rent']].mean()
3 years2.columns = ['Year', 'Rent']
4
5 # Create a Line plot
6 xvals = ['2021','2022','2023','2024']
7 ax = sns.lineplot(data=years2, x=xvals, y="Rent", color = 'darkred')
8 plt.title("Average Rents by Year")
9 ax.spines.right.set_visible(False)
10 ax.spines.top.set_visible(False);
```



Average rents peaked in 2022 where number of listings is low as we might expect. Since then it has declined as listings increase. However, the rates of change (decrease in average rents and increase in properties listed) differ, with average rents declining less rapidly than the number of listings. This might suggest continuing strong demand for rental properties not matched by supply. We can look at the number of days listings over time later which might shed some more light on this.

```
In [540]: 1 # Average rents by Month
2 years4 = dubai.groupby('MonthName',as_index = False)[['Rent']].mean().sort_values(by = 'Rent',ascending =
3
4 # Create a Line plot
5 ax = sns.barplot(data=years4, y='MonthName', x="Rent",palette = 'coolwarm_r')
6 plt.title("Average Rents by Month")
7 ax.spines.right.set_visible(False)
8 ax.spines.top.set_visible(False);
```



Average rental prices are high in August and June where the number of rentals being listed are low and demand is strong and low in March and April where there are more listings. The exception is May where the number of listings is low and the rental values are also lowest reflecting lower demand for properties and lower supply on the market.

### 3.4 Age of Listing

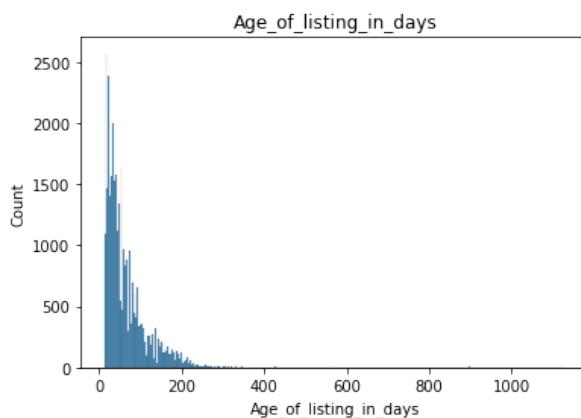
```
In [541]: 1 # Basic stats for age of listing  
2 df2['Age_of_listing_in_days'].describe()
```

```
Out[541]: count    34229.000000  
mean      63.917789  
std       55.183479  
min       12.000000  
25%      27.000000  
50%      46.000000  
75%      81.000000  
max     1131.000000  
Name: Age_of_listing_in_days, dtype: float64
```

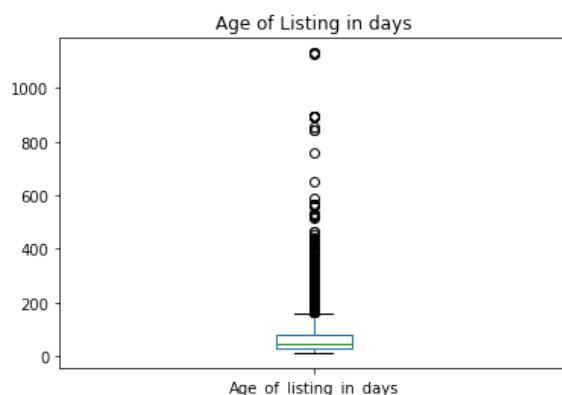
```
In [542]: 1 # Skew  
2 df2['Age_of_listing_in_days'].skew()
```

```
Out[542]: 3.7106414442965345
```

```
In [543]: 1 # Plot the distribution of area in square feet  
2 sns.histplot(data = df2['Age_of_listing_in_days'])  
3 plt.title('Age_of_listing_in_days');
```



```
In [544]: 1 # Boxplot of the data  
2 df2['Age_of_listing_in_days'].plot(kind = 'box')  
3 plt.title('Age of Listing in days');
```



The longest listing is for 1131 days, the shortest for just 12. The median listing time is 46 days and mean of 64 days which is a couple of months. Again the data is right skewed with the median values below the mean.

```
In [545]: 1 # Listings over 1000 days
2 df2[df2['Age_of_listing_in_days']>1000]
```

Out[545]:

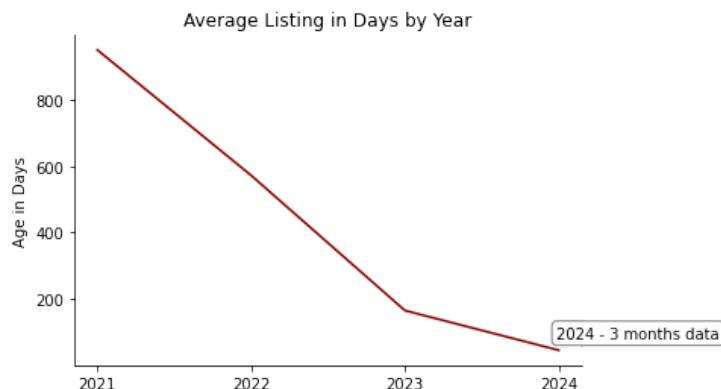
|       | Address   | Rent  | Beds | Baths | Type      | Area_in_sqft | Rent_per_sqft | Rent_category | Furnishing  | Posted_date | Age_of_list |
|-------|---|-------|------|-------|-----------|--------------|---------------|---------------|-------------|-------------|-------------|
| 32580 | CBD-F05, Central Business District, Internatio... | 53000 | 2    | 2     | Apartment | 915          | 57.923497     | Low           | Unfurnished | 2021-03-17  |             |
| 32569 | B-12, China Cluster, International City, Dubai    | 53000 | 2    | 2     | Apartment | 915          | 57.923497     | Low           | Unfurnished | 2021-03-17  |             |
| 34391 | A-09, China Cluster, International City, Dubai    | 53000 | 2    | 2     | Apartment | 915          | 57.923497     | Low           | Unfurnished | 2021-03-18  |             |
| 32590 | P-01, France Cluster, International City, Dubai   | 31500 | 0    | 2     | Apartment | 484          | 65.082645     | Low           | Unfurnished | 2021-03-17  |             |

These all seem to be low rent apartments in one location at an address of the China Cluster or France Cluster. There is clearly something about these rentals being in this location that is important so it wouldnt be good idea to just remove them.

```
In [546]: 1 # Change over time?
2 years5 = dubai.groupby('Year',as_index = False)[['Age_of_listing_in_days']].mean()
3 years5.columns = ['Year','Age in Days']
4 display(years5)
5
6 # Create a Line plot
7 xvals = ['2021','2022','2023','2024']
8 ax = sns.lineplot(data=years5, x=xvals, y="Age in Days", color = 'darkred')
9 plt.title("Average Listing in Days by Year")
10 ax.spines.right.set_visible(False)
11 ax.spines.top.set_visible(False)
12 props = dict(boxstyle='round', facecolor='White', alpha=0.5)
13 ax.text(0.95, 0.12, '2024 - 3 months data', transform=ax.transAxes, fontsize=10,
14         verticalalignment='top', bbox=props)
15 ;
```

|   | Year | Age in Days |
|---|------|-------------|
| 0 | 2021 | 948.933333  |
| 1 | 2022 | 571.800000  |
| 2 | 2023 | 165.659463  |
| 3 | 2024 | 46.344788   |

Out[546]: ''

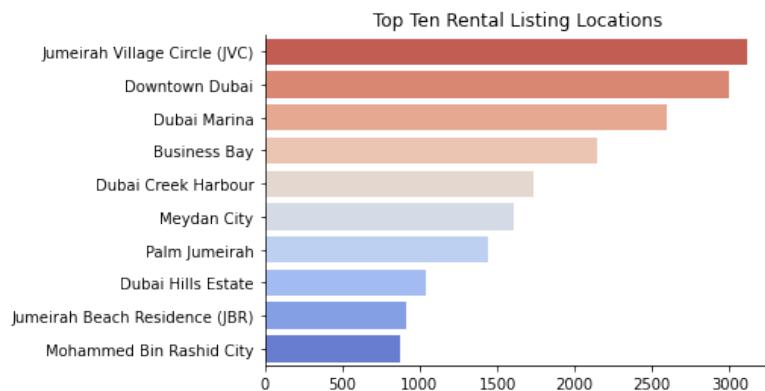


We can see that average listing times by year have fallen over the period from an average of 949 days in 2021 (9 month average) to 166 days in 2023. The rate of decline slows between 2023 and 2024 but again, and not as sharp as for the rise in number of listings. This suggests that although many more properties are being listed, that demand is still very strong especially into 2024 where

properties are on the market for around 6 weeks on average.

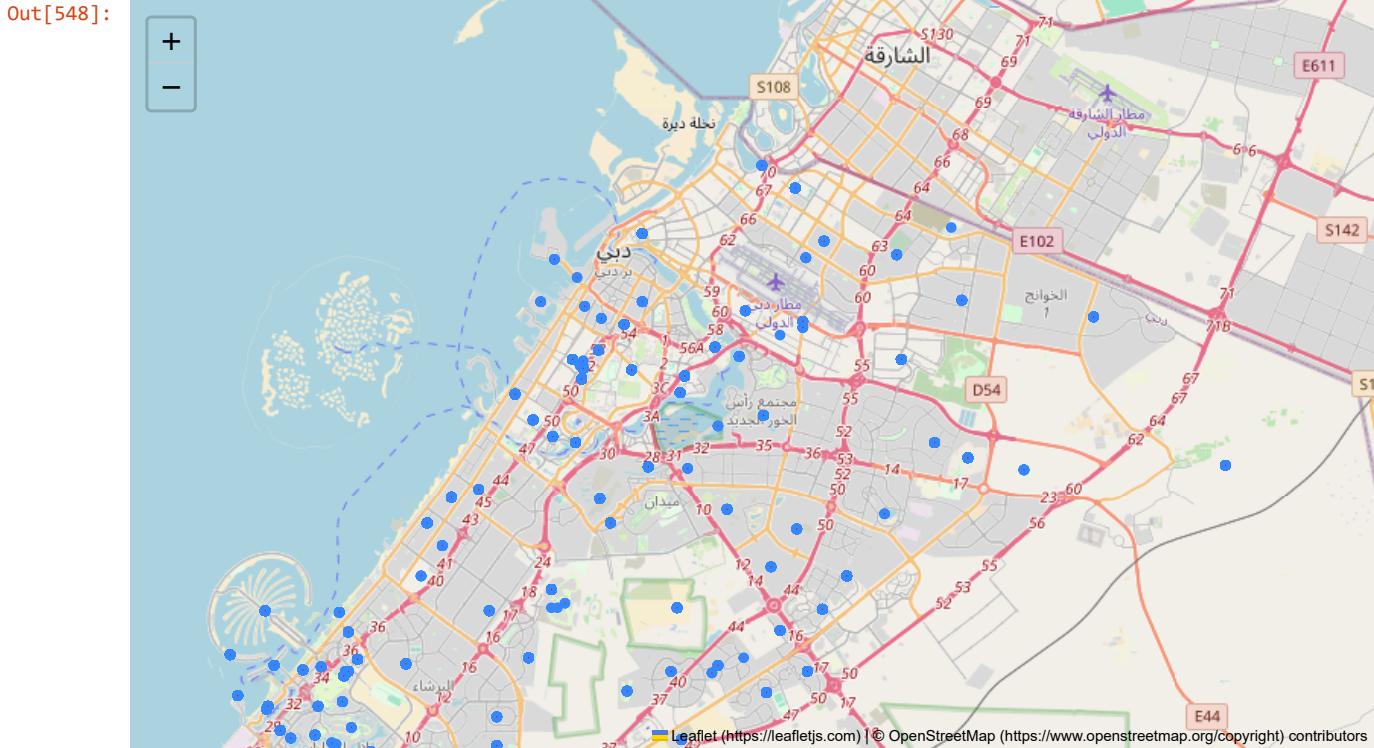
### 3.5 Location

```
In [547]: 1 # Top ten Locations being listed
2 top_10_locations = pd.DataFrame(dubai['Location'].value_counts().head(10).reset_index())
3
4 # Plot Top ten Locations
5 ax = sns.barplot(data = top_10_locations, y = 'index', x = 'Location', palette = 'coolwarm_r')
6 plt.title('Top Ten Rental Listing Locations')
7 plt.ylabel(None)
8 plt.xlabel(None)
9 ax.spines.right.set_visible(False)
10 ax.spines.top.set_visible(False);
```



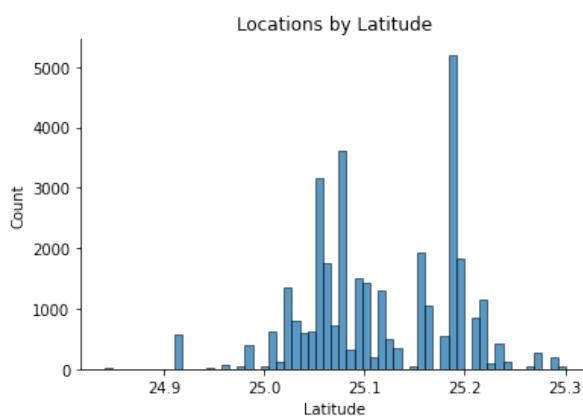
Most rentals listed for the period are for Jumeirah Village, Downtown Dubai and the Dubai Marina.

```
In [548]: 1 # Create a map of Locations for rentals in Dubai
2
3 # Create a base map
4 m = folium.Map(location = [dubai['Latitude'].mean(), dubai['Longitude'].mean()],zoom_start = 10)
5
6 # Add points to the map
7 for idx, row in dubai.iterrows():
8     folium.CircleMarker(location = (row['Latitude'], row['Longitude']),
9                          radius = 2,
10                         weight = 1,
11                         fill = True,
12                         fill_color = 'red',
13                         fill_opacity = 0.7).add_to(m)
14
m
```

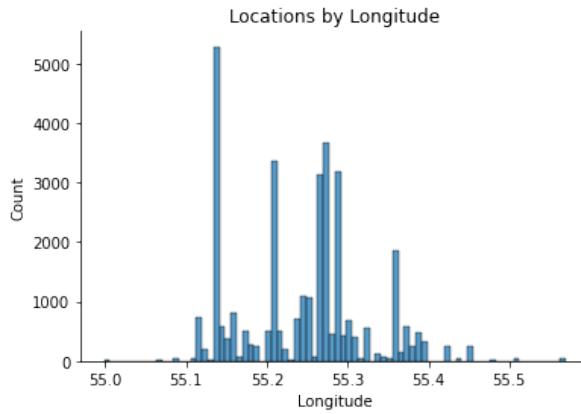


The map shows there are two clusters of properties which we can see if we plot the latitudes and longitudes showing the north south and east west clusters

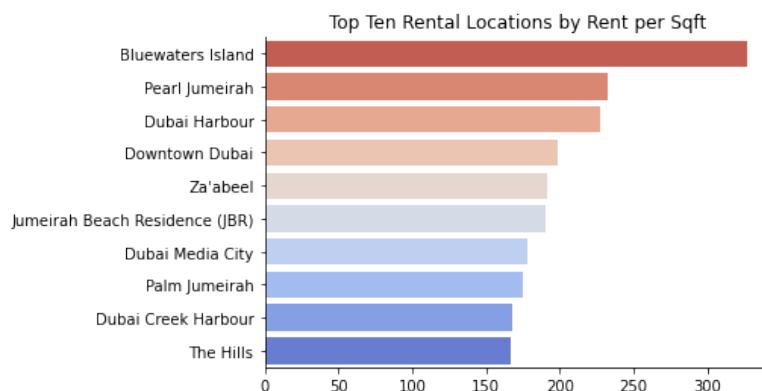
```
In [549]: 1 # We can plot the Latitudes
2 ax=sns.histplot(dubai, x = 'Latitude')
3 plt.title('Locations by Latitude')
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



```
In [550]: 1 # We can plot the Longitudes
2 ax=sns.histplot(dubai, x = 'Longitude')
3 plt.title('Locations by Longitude')
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```



```
In [551]: 1 # What are the top ten Locations by rent per square foot?
2 top_10_locations_rent = dubai.groupby('Location',as_index = False)[['Rent_per_sqft']].mean().sort_values
3
4 # Plot Top ten locations
5 ax = sns.barplot(data = top_10_locations_rent, y = 'Location', x = 'Rent_per_sqft', palette = 'coolwarm')
6 plt.title('Top Ten Rental Locations by Rent per Sqft')
7 plt.ylabel(None)
8 plt.xlabel(None)
9 ax.spines.right.set_visible(False)
10 ax.spines.top.set_visible(False);
```



```
In [552]: 1 dubai[dubai['Location']=='Bluewaters Island']['Rent_per_sqft'].mean()
```

Out[552]: 327.33516697397675

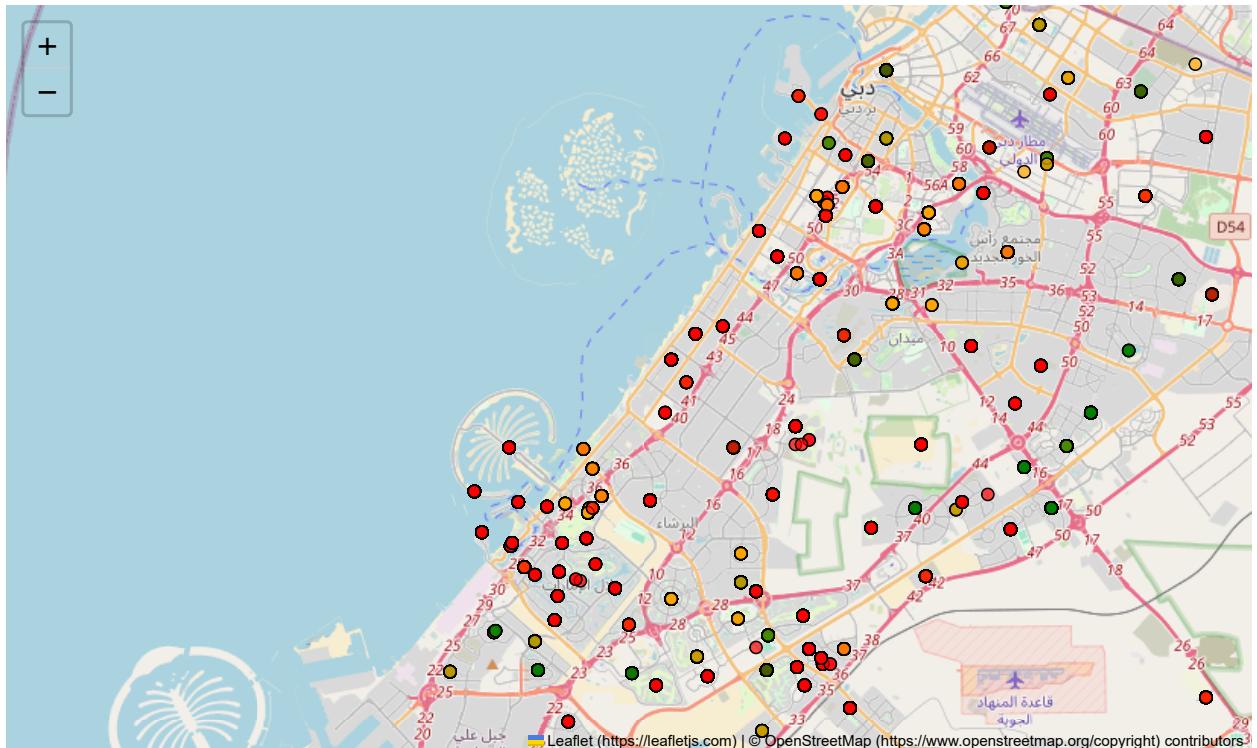
We can see that Bluewaters Island has the highest average rent per square ft of 329 AED per square ft.

```

In [997]: 1 # Lets add a map that shows locations but categorised by the rent category or high, medium or low
2
3 # Create a base map
4 m = folium.Map(location = [dubai['Latitude'].mean(), dubai['Longitude'].mean()],zoom_start = 11)
5
6 # Define categories colour mapping
7 color_map = {'High':'red', 'Medium':'orange', 'Low':'green'}
8
9 # Add points to the map
10 for idx, row in dubai.iterrows():
11     folium.CircleMarker(location = (row['Latitude'], row['Longitude']),
12                         radius = 4,
13                         color = 'black',
14                         weight = 1,
15                         color_map = color_map[row['Rent_category']],
16                         fill = True,
17                         fill_color = color_map[row['Rent_category']],
18                         fill_opacity = 0.7).add_to(m)
19
m

```

Out[997]:



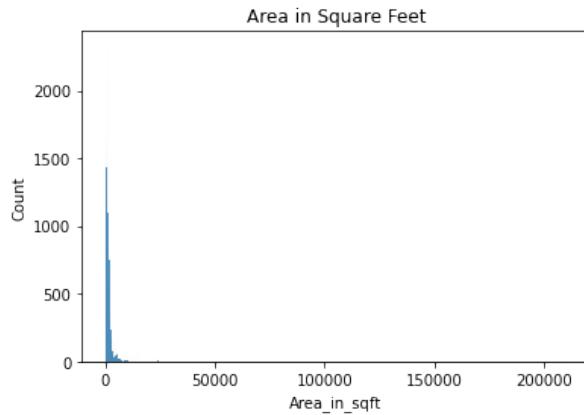
The map shows that most of the rentals are high priced and that they are clustered along the water and to the south of the central part of the city. Lower priced rentals are a bit further out, although there are some exceptions. Medium priced rentals are scattered about the city with clusters in the centre north and centre south.

```
In [554]: 1 dubai.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 33895 entries, 29068 to 63317
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Address          33895 non-null   object  
 1   Rent              33895 non-null   int64   
 2   Beds              33895 non-null   int64   
 3   Baths             33895 non-null   int64   
 4   Type              33895 non-null   object  
 5   Area_in_sqft     33895 non-null   int64   
 6   Rent_per_sqft    33895 non-null   float64 
 7   Rent_category    33895 non-null   object  
 8   Furnishing       33895 non-null   object  
 9   Posted_date       33895 non-null   datetime64[ns]
 10  Age_of_listing_in_days 33895 non-null   int64   
 11  Location          33895 non-null   object  
 12  Latitude          33895 non-null   float64 
 13  Longitude         33895 non-null   float64 
 14  Year              33895 non-null   int64   
 15  Month             33895 non-null   int64   
 16  MonthName         33895 non-null   object  
 17  RentZScore        33895 non-null   float64 
dtypes: datetime64[ns](1), float64(4), int64(7), object(6)
memory usage: 5.9+ MB
```

### 3.6 Area

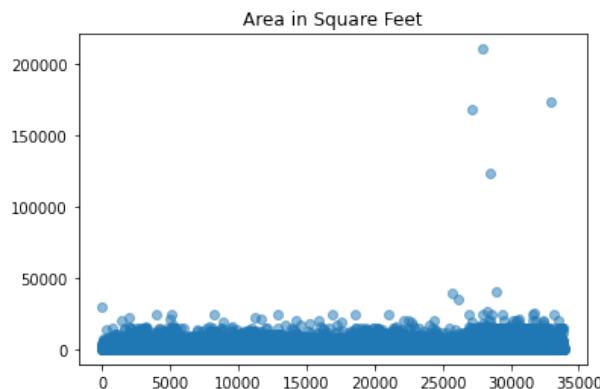
```
In [555]: 1 # Plot the distribution of area in square feet
2 sns.histplot(data = dubai['Area_in_sqft'])
3 plt.title('Area in Square Feet');
```



We have a lot of right skew in the data, indicating some very large properties. We can have a look what these are

#### Outlier Investigation

```
In [557]: 1 # Scatter plot of the property areas
2 ind = list(range(0,33895))
3 vals = list(dubai['Area_in_sqft'])
4
5 fig, ax = plt.subplots(figsize = (6,4))
6 plt.scatter(ind,vals, alpha = 0.5)
7
8 plt.title('Area in Square Feet');
```



Several properties are above 100000 square feet

```
In [558]: 1 # Split into those over and those under or equal to 100000 square feet
2 over = dubai[dubai['Area_in_sqft']>=100000].sort_values(by = 'Area_in_sqft', ascending = False)
3 under = dubai[dubai['Area_in_sqft']<100000].sort_values(by = 'Area_in_sqft', ascending = False)
```

```
In [559]: 1 over
```

Out[559]:

|  |       | Address   | Rent    | Beds | Baths | Type             | Area_in_sqft | Rent_per_sqft | Rent_category | Furnishing  | Posted_date | Age_of_property |
|--|-------|---|---------|------|-------|------------------|--------------|---------------|---------------|-------------|-------------|-----------------|
|  | 57236 | Marina Residences<br>1, Marina Residences,<br>Palm J... | 900000  | 4    | 2     | Penthouse        | 210254       | 4.280537      | High          | Unfurnished | 2024-03-20  |                 |
|  | 62394 | One Za'abeel The Residences,<br>One Za'abeel,<br>Za'... | 1200000 | 1    | 1     | Residential Plot | 173565       | 6.913836      | High          | Unfurnished | 2024-03-11  |                 |
|  | 56445 | Jumeirah Village Triangle (JVT),<br>Dubai               | 299999  | 6    | 2     | Villa            | 168046       | 1.785220      | High          | Furnished   | 2024-04-03  |                 |
|  | 57810 | Jumeirah 3,<br>Jumeirah,<br>Dubai                       | 420000  | 5    | 2     | Villa Compound   | 123696       | 3.395421      | High          | Unfurnished | 2024-03-15  |                 |



```
In [560]: 1 # Create list of these Locations where property is over 100000
2 locations = over['Location'].tolist()
3 print(locations)
4
5 # Get the mean areas for these Locations for the rentals excluding these properties
6 display(under[under['Location'].isin(locations)].groupby('Location')['Area_in_sqft'].mean())
7
8 # Get the median areas for these Locations for the rentals excluding these properties
9 display(under[under['Location'].isin(locations)].groupby('Location')['Area_in_sqft'].median())
```

[ 'Palm Jumeirah', "Za'abeel", 'Jumeirah Village Triangle (JVT)', 'Jumeirah' ]

| Location                        |             |
|---------------------------------|-------------|
| Jumeirah                        | 3735.671388 |
| Jumeirah Village Triangle (JVT) | 2587.018349 |
| Palm Jumeirah                   | 1960.056367 |
| Za'abeel                        | 1319.143820 |

Name: Area\_in\_sqft, dtype: float64

| Location                        |      |
|---------------------------------|------|
| Jumeirah                        | 2907 |
| Jumeirah Village Triangle (JVT) | 1850 |
| Palm Jumeirah                   | 1582 |
| Za'abeel                        | 1200 |

Name: Area\_in\_sqft, dtype: int64

Comparing these large properties to the mean and median values for every other property in those respective neighbourhoods, it seems that they are very much outliers in those areas so we will drop them from the data.

```
In [737]: 1 # Create a copy of the revised data
2 df2 = under.copy()
3 df2.shape
```

Out[737]: (33891, 18)

## 3.7 Relationships

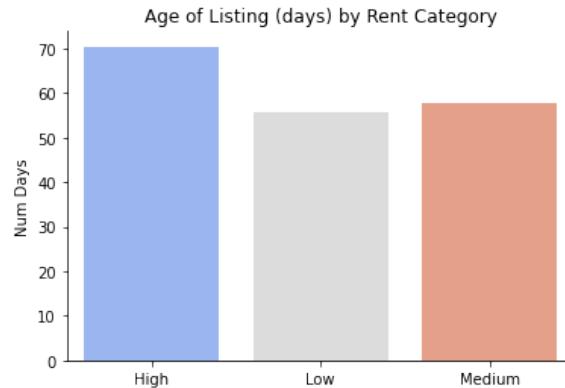
### Listing Age and Rents

```
In [738]: 1 # What category of rentals is on the market for the longest and which the shortest?
2
3 rent_age = df2.groupby('Rent_category',as_index = False)[['Age_of_listing_in_days']].mean()
4 rent_age
```

Out[738]:

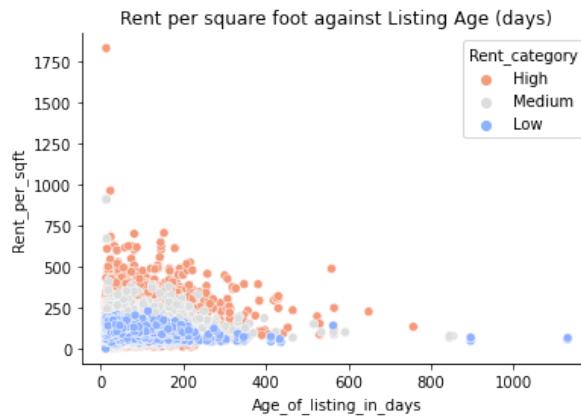
| Rent_category | Age_of_listing_in_days |
|---------------|------------------------|
| 0             | 70.428125              |
| 1             | 55.581395              |
| 2             | 57.676539              |

```
In [739]: 1 # Plot the data
2 ax = sns.barplot(data = rent_age, x = 'Rent_category', y = 'Age_of_listing_in_days', palette = 'coolwarm'
3 plt.title("Age of Listing (days) by Rent Category")
4 plt.xlabel(None)
5 plt.ylabel('Num Days')
6 ax.spines.right.set_visible(False)
7 ax.spines.top.set_visible(False);
```



Listing days gives us a sense of demand in the market. High rental properties are on the market an average of 70 days compared to medium priced rentals at 58 and low priced at 58. Therefore there is very little difference between the low priced and medium priced categories in terms of the number of days they are on the market. Our earlier chart showed about double the number of listings for medium priced rentals as low priced and yet they remain on the market for a similar time, suggesting the demand might be primarily for the medium priced rental properties.

```
In [740]: 1 # Listing age against rent per square foot scatterplot
2 ax = sns.scatterplot(data = df2, x = 'Age_of_listing_in_days', y = 'Rent_per_sqft', hue = 'Rent_category'
3 plt.title('Rent per square foot against Listing Age (days)')
4 ax.spines.right.set_visible(False)
5 ax.spines.top.set_visible(False);
```

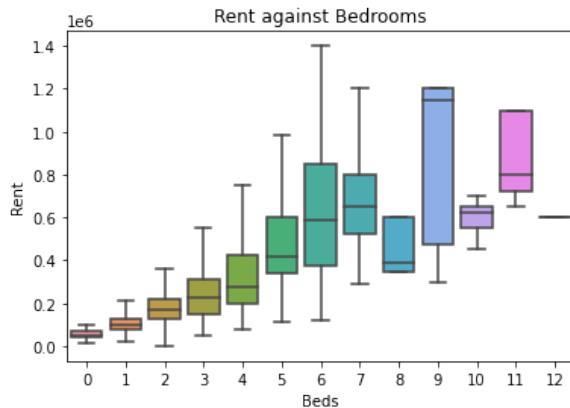


We can see that there is some weak positive correlation here. There are outliers in the data with low rental property on the market for more than 1000 days and higher priced appearing to be snapped up immediately. This should probably be the focus of additional work to identify what issues if any relate to their being empty for so long as clearly there are issues other than price affecting the rentability for some of the lower priced ones.

### Other Relationships to Rent

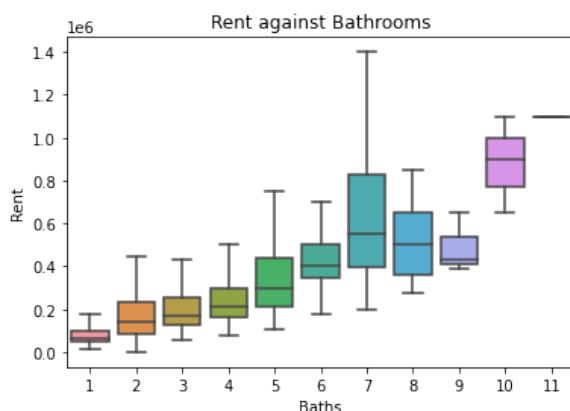
The feature we would want to predict would be the rent. We want to identify those other features that have any kind of correlation with rent. For example does the rent value have any correlation with locational data, or with the number of bedrooms?

```
In [741]: 1 # Bedrooms
2 sns.boxplot(data = df2, x = 'Beds', y = 'Rent', showfliers = False)
3 plt.title("Rent against Bedrooms");
```



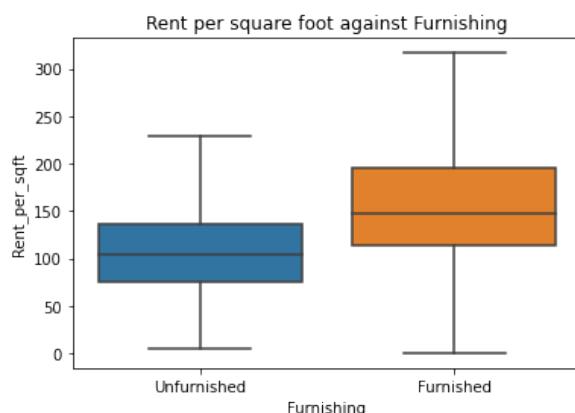
There appears to be a positive correlation between number of beds and average rental values up to about 7, whereby more bedrooms attracts higher rents. Of interest is the wider range of rental values for properties with large numbers of bedrooms (above 5) from higher rental top-end through to lower rental.

```
In [742]: 1 # Bathrooms
2 sns.boxplot(data = df2, x = 'Baths', y = 'Rent', showfliers = False)
3 plt.title("Rent against Bathrooms");
```



More bathrooms appear to attract higher rents, up to a point. At 7 bedrooms there is again a wider range of rental values

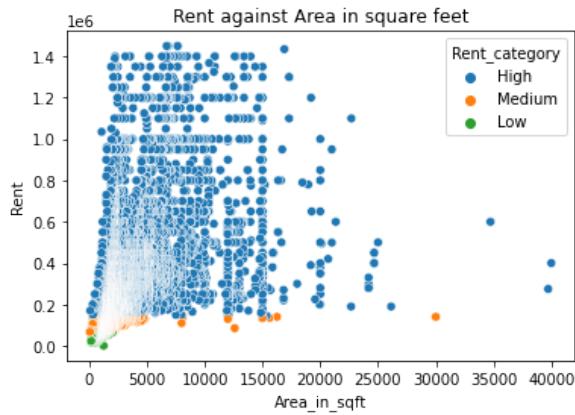
```
In [743]: 1 # Compare furnishing with rental values excluding outliers
2 sns.boxplot(data = df2, x = 'Furnishing', y = 'Rent_per_sqft', showfliers = False)
3 plt.title("Rent per square foot against Furnishing");
```



Furnished properties have a higher rental value than unfurnished

In [744]:

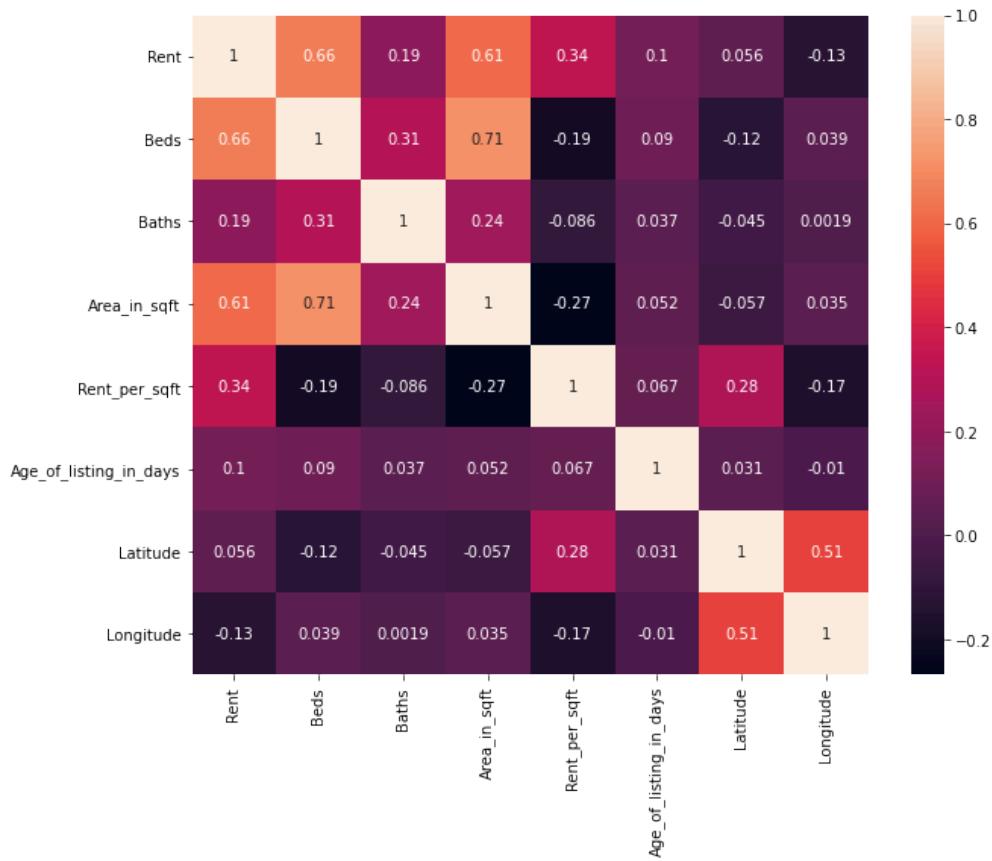
```
1 # Area against rent
2 sns.scatterplot(data = df2, x = 'Area_in_sqft', y = 'Rent', hue = 'Rent_category')
3 plt.title('Rent against Area in square feet');
```



As we might expect there seems to be a positive correlation between the area and rents

In [745]:

```
1 # Correlations in a heatmap
2 cols = ['Rent', 'Beds', 'Baths', 'Area_in_sqft', 'Rent_per_sqft', 'Age_of_listing_in_days', 'Latitude', 'Longitude']
3 fig, ax = plt.subplots(figsize = (10,8))
4 ax = sns.heatmap(df2[cols].corr(), annot = True)
```



Rent is strongly positively correlated with area and number of bedrooms, negatively correlated with location (longitude), so that moving inland has lower rents and weakly positively correlated with the number of bathrooms and the age of listing. Location north or south indicated by latitude has a very small positive correlation.

From the analysis, we might look at including most of these features but exclude the geo-coordinates

## 4.0 Prediction

### 4.1 Data Preparation

```
In [883]: 1 df3 = df2.copy()
```

Trying to avoid a 'kitchen sink' regression, we will pick some of the features to include. The address has many different values which would complicate the analysis and the information is available in other features. The rent per square foot is contains information already for the target of rent we are trying to predict. The posted date is effectively covered by the year and month and the location will be covered by the latitude and longitude. The monthname copies the month feature and the rentzscore is added earlier and is irrelevant.

```
In [884]: 1 # Drop columns we don't need
2 df3.drop(columns = ['Address', 'Rent_per_sqft', 'Posted_date', 'Location', 'MonthName', 'RentZScore'], axis =
```

```
In [885]: 1 # Reset the index so we can manipulate the dataframe
2 df3.reset_index(inplace = True, drop = True)
3
4 # Split dataframe into X and y (features and target)
5 X = df3.drop(columns = ['Rent'], axis = 1)
6 y = df3[['Rent']]
```

We have both ordinal and nominal categorical values as well as numerical. For example the rent category is ordered from low to high, however the furnishing and type are NOT ordered so we would not want to apply a sequential ranking to these so would choose a one-hot encoding.

```
In [886]: 1 # Replace categorical ordinal values in rent category with numerical
2 X['Rent_category'].replace({'High':3, 'Medium':2, 'Low':1}, inplace = True)
3
4 # One hot encoding for the type, furnishing and month (we set drop first = true to avoid multicollinear
5 X = pd.get_dummies(X, prefix=None, columns=['Type', 'Furnishing', 'Year', 'Month'], drop_first = True)
```

We will standardise the numerical features but as we will fit and transform on the training data and then should only transform on the test data, we will need to split the dataframe into training and test sets first

```
In [887]: 1 # Set up the training and test sets
2 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size = 0.2)
3 print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

```
(27112, 30) (6779, 30) (27112, 1) (6779, 1)
```

```
In [888]: 1 # Split off the numerical features and standardise these
2 nums = ['Beds', 'Baths', 'Area_in_sqft', 'Age_of_listing_in_days', 'Latitude', 'Longitude']
3 X_train_num = X_train[nums]
4 X_test_num = X_test[nums]
5
6 # Standardise the numeric features (Beds, Bath, Area, Age of Listing), fit transform the training set
7 scale = StandardScaler()
8 X_train_num = pd.DataFrame(scale.fit_transform(X_train_num))
9 X_train_num.columns = nums
10
11 # Transform the test set
12 X_test_num = pd.DataFrame(scale.transform(X_test_num))
```

```
In [889]: 1 # Reset the index on each dataframe
2 X_train.reset_index(drop=True, inplace=True)
3 X_train_num.reset_index(drop=True, inplace=True)
4
5 X_test.reset_index(drop=True, inplace=True)
6 X_test_num.reset_index(drop=True, inplace=True)
7
8 # Drop the numerical columns in the original dataframe
9 X_train.drop(columns = nums, axis = 1, inplace = True)
10 X_test.drop(columns = nums, axis = 1, inplace = True)
11
12 # Concat the scaled data back to the original dataframe
13 X_train = pd.concat([X_train, X_train_num], axis = 1)
14 X_test = pd.concat([X_test, X_test_num], axis = 1)
```

We now have our training and test sets ready for modelling

## 4.2 Multiple Linear Regression

```
In [890]: 1 # Fit the Linear model to the training data and get the training score
2 lr = LinearRegression()
3 mod1 = lr.fit(X_train, y_train)
4 mod1.score(X_train, y_train)
```

Out[890]: 0.5874595003889362

```
In [891]: 1 # Predict on the test set
2 y_pred = mod1.predict(X_test)
3 r2_score(y_test, y_pred)
```

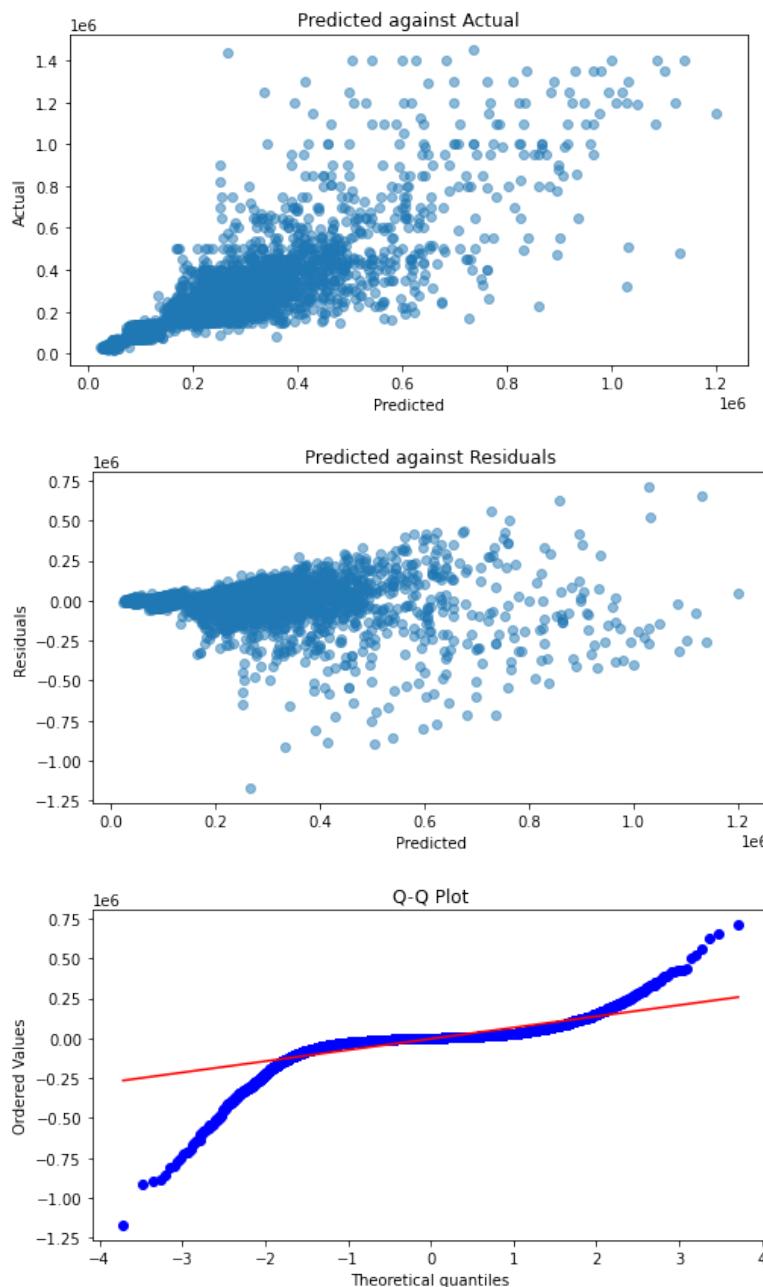
Out[891]: 0.6015289766052676

From this first model, 61% of the variation in rental values can be explained by our model. There are clearly other factors at play in the market affecting the price. We should look at the residuals as we have made some assumptions using a linear regression on this data

```
In [934]: 1 # Concat the actual and predicted and calculate residuals
2 results = pd.concat([y_test, y_pred],axis = 1)
3 results.columns = ['Actual','Predicted']
4 results['Res'] = results['Predicted'] - results['Actual']
```

In [948]:

```
1 # Plot predicted against actual
2 fig, ax = plt.subplots(figsize = (8,4))
3 plt.scatter(results['Predicted'], results['Actual'], alpha = 0.5)
4 plt.xlabel('Predicted')
5 plt.ylabel('Actual')
6 plt.title('Predicted against Actual');
7
8 # Predicted against residual
9 fig, ax = plt.subplots(figsize = (8,4))
10 plt.scatter(results['Predicted'], results['Res'], alpha = 0.5)
11 plt.xlabel('Predicted')
12 plt.ylabel('Residuals')
13 plt.title('Predicted against Residuals');
14
15 # QQ Plot
16 plt.figure(figsize=(8, 4))
17 stats.probplot(results['Res'], dist="norm", plot=plt)
18 plt.title('Q-Q Plot');
19
```



We can see from these plots that there is increasing variance in the residuals indicating heteroscedasticity with the residuals getting larger across the plot. This reflects the fact that rents can be high or low depending on the other features so we cannot assume they are high or low. The QQ plot shows the blue line varying from the normal red line at both ends. This goes against one of the assumptions of the linear model, so we can either try to adjust this model with maybe a log transform or try different models that do not require these assumptions. We will look at other models.

## 4.3 Other models

We will look at a couple of other models accepting the default parameters.

### Nearest Neighbours

```
In [898]: 1 # Fit model and get training score  
2 knn = neighbors.KNeighborsRegressor(n_neighbors = 5)  
3 mod2 = knn.fit(X_train, y_train)  
4 mod2.score(X_train, y_train)
```

Out[898]: 0.8286433707539533

```
In [899]: 1 # Predict on the test set  
2 y_pred = mod2.predict(X_test)  
3 r2_score(y_test, y_pred)
```

Out[899]: 0.7517239157864216

This model seems to have performed better than our linear regression model

### Random Forest

```
In [952]: 1 clf = RandomForestClassifier()  
2 mod3 = clf.fit(X_train, y_train)  
3 mod3.score(X_train, y_train)
```

Out[952]: 0.9806727648273827

```
In [953]: 1 # Predict on the test set  
2 y_pred = mod3.predict(X_test)  
3 r2_score(y_test, y_pred)
```

Out[953]: 0.7349611182082432

The higher training score compared to test score suggests this model might be overfitting.

```
In [965]: 1 # Reduce number of estimators - various values were explored here to see if we could get a better test  
2 clf = RandomForestClassifier(n_estimators = 2)  
3 mod3 = clf.fit(X_train, y_train)  
4 print('Training Score:', mod3.score(X_train, y_train))  
5  
6 # Predict on the test set  
7 y_pred = mod3.predict(X_test)  
8 r2_score(y_test, y_pred)
```

Training Score: 0.664281498967247

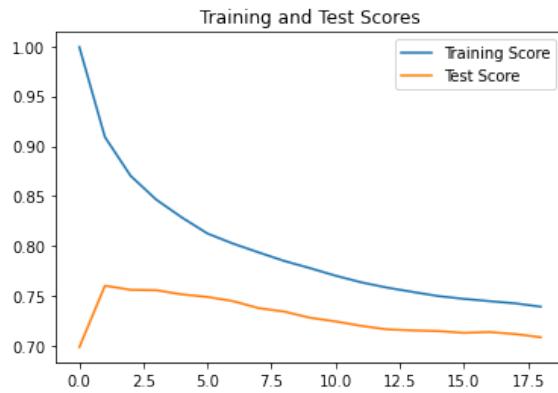
Out[965]: 0.6621000484924996

We have reduced overfitting but the training and test scores are not as good as for the nearest neighbours algorithm. We can see if we can improve the results by trying different numbers of neighbours for that model.

### Nearest Neighbors - parameter tuning

```
In [977]: 1 # Fit model and get training score
2 n_neighbors = list(range(1,20,1))
3 trains = []
4 tests = []
5
6 for i in n_neighbors:
7     knn = neighbors.KNeighborsRegressor(n_neighbors = i)
8     mod2 = knn.fit(X_train, y_train)
9     y_pred = mod2.predict(X_test)
10
11    train_score = mod2.score(X_train, y_train)
12    test_score = r2_score(y_test, y_pred)
13    trains.append(train_score)
14    tests.append(test_score)
```

```
In [984]: 1 plt.plot(trains)
2 plt.plot(tests)
3 plt.title('Training and Test Scores')
4 plt.legend(['Training Score', 'Test Score']);
```



```
In [993]: 1 scores = pd.DataFrame(list(zip(trains, tests)))
2 scores['Diff'] = scores[0] - scores[1]
3 scores.sort_values(by = 'Diff', ascending = False)
```

Out[993]:

|    | 0    | 1    | Diff |
|----|------|------|------|
| 0  | 1.00 | 0.70 | 0.30 |
| 1  | 0.91 | 0.76 | 0.15 |
| 2  | 0.87 | 0.76 | 0.11 |
| 3  | 0.85 | 0.76 | 0.09 |
| 4  | 0.83 | 0.75 | 0.08 |
| 5  | 0.81 | 0.75 | 0.06 |
| 6  | 0.80 | 0.74 | 0.06 |
| 7  | 0.79 | 0.74 | 0.06 |
| 8  | 0.79 | 0.73 | 0.05 |
| 9  | 0.78 | 0.73 | 0.05 |
| 10 | 0.77 | 0.72 | 0.05 |
| 11 | 0.76 | 0.72 | 0.04 |
| 12 | 0.76 | 0.72 | 0.04 |
| 13 | 0.75 | 0.72 | 0.04 |
| 14 | 0.75 | 0.72 | 0.03 |
| 15 | 0.75 | 0.71 | 0.03 |
| 16 | 0.74 | 0.71 | 0.03 |
| 17 | 0.74 | 0.71 | 0.03 |
| 18 | 0.74 | 0.71 | 0.03 |

Looking at the training and test scores there is no further improvement beyond 14 neighbours

In [994]:

```
1 # Fit model and get training score
2 knn = neighbors.KNeighborsRegressor(n_neighbors = 14)
3 mod2 = knn.fit(X_train, y_train)
4 print(mod2.score(X_train, y_train))
5
6 # Predict on the test set
7 y_pred = mod2.predict(X_test)
8 print(r2_score(y_test, y_pred))
```

```
0.7542007165680242
0.715708109253427
```

There is plenty more we could do with this project but for now we can conclude that with the K Nearest Neighbours Algorithm that we can explain about 72% of the variation in rental prices from the model including the features we have chosen. This leaves another 28% explainable by other factors we have not considered or do not have access to. This might include all kinds of other things such as location close to transport routes, schools, amenities, local conditions, market conditions and government policies. This shows how complicated building a good model to predict rental prices is.