

Name: Raj C. Mhatre

Email id :rmhatre404@gmail.com (mailto:rmhatre404@gmail.com) , contact No.: +91-9768877272

Linkedin ID: <https://www.linkedin.com/in/raj-c-mhatre-571b67219/>

▼ Data Science Regression Project: Predicting Home Prices in Bangalore

Dataset is downloaded from here: <https://www.kaggle.com/amitabhajoy/bengaluru-house-price-data>

Importing all the required libraries:

```
import warnings
warnings.filterwarnings("ignore")

import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import matplotlib
matplotlib.rcParams["figure.figsize"] = (20,10)
```

Data Load: Load banglore home prices into a dataframe

```
df1 = pd.read_csv("/content/Bengaluru_House_Data.csv")
df1
```

	area_type	availability	location	size	society	total_sqft	bath	balcony	price
0	Super built-up Area	19-Dec	Electronic City Phase II	2 BHK	Coomee	1056	2.0	1.0	39.07
1	Plot Area	Ready To Move	Chikka Tirupathi	4 Bedroom	Theanmp	2600	5.0	3.0	120.00
2	Built-up Area	Ready To Move	Uttarahalli	3 BHK	NaN	1440	2.0	3.0	62.00
3	Super built-up Area	Ready To Move	Lingadheeranahalli	3 BHK	Soiewre	1521	3.0	1.0	95.00
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	NaN	1200	2.0	1.0	51.00
...
13315	Built-up Area	Ready To Move	Whitefield	5 Bedroom	ArsiaEx	3453	4.0	0.0	231.00
13316	Super built-up Area	Ready To Move	Richards Town	4 BHK	NaN	3600	5.0	NaN	400.00
...

```
df1.shape
```

```
(13320, 9)
```

```
df1["area_type"].value_counts()
```

```
Super built-up Area    8790
Built-up Area          2418
Plot Area              2025
Carpet Area             87
Name: area_type, dtype: int64
```

Drop features that are not required to build our model

```
df2 = df1.drop(['area_type', 'society', 'balcony', 'availability'], axis='columns')
df2.head()
```

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Tirupathi	4 Bedroom	2600	5.0	120.00
2	Uttarahalli	3 BHK	1440	2.0	62.00
3	Lingadheeranahalli	3 BHK	1521	3.0	95.00



▼ Data Cleaning: Handle NA values

```
df2.isnull().sum()
```

```
location      1
size          16
total_sqft    0
bath          73
price         0
dtype: int64
```

```
df3 = df2.dropna()
```

```
df3.head()
```

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Tirupathi	4 Bedroom	2600	5.0	120.00
2	Uttarahalli	3 BHK	1440	2.0	62.00
3	Lingadheeranahalli	3 BHK	1521	3.0	95.00
4	Kothanur	2 BHK	1200	2.0	51.00



```
df3.isnull().sum()
```

```
location      0
size          0
total_sqft    0
bath          0
price         0
dtype: int64
```

```
df3['size'].unique()
```

```
array(['2 BHK', '4 Bedroom', '3 BHK', '4 BHK', '6 Bedroom', '3 Bedroom',
       '1 BHK', '1 RK', '1 Bedroom', '8 Bedroom', '2 Bedroom',
       '7 Bedroom', '5 BHK', '7 BHK', '6 BHK', '5 Bedroom', '11 BHK',
       '9 BHK', '9 Bedroom', '27 BHK', '10 Bedroom', '11 Bedroom',
       '10 BHK', '19 BHK', '16 BHK', '43 Bedroom', '14 BHK', '8 BHK',
       '12 Bedroom', '13 BHK', '18 Bedroom'], dtype=object)
```

```
df3['bhk'] = df3['size'].apply(lambda x: int(x.split(' ')[0]))
```

```
df3.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600	5.0	120.00	4
2	Uttarahalli	3 BHK	1440	2.0	62.00	3
3	Lingadheeranahalli	3 BHK	1521	3.0	95.00	3
4	Kothanur	2 BHK	1200	2.0	51.00	2



▼ Feature Engineering

Add new feature(integer) for bhk (Bedrooms Hall Kitchen)

```
df3['bhk'].unique()

array([ 2,  4,  3,  6,  1,  8,  7,  5, 11,  9, 27, 10, 19, 16, 43, 14, 12,
        13, 18])

df3[df3.bhk>20]
```

	location	size	total_sqft	bath	price	bhk
1718	2Electronic City Phase II	27 BHK	8000	27.0	230.0	27
4684	Munnekollal	43 Bedroom	2400	40.0	660.0	43

```
df3.total_sqft.unique()

array(['1056', '2600', '1440', ..., '1133 - 1384', '774', '4689'],
      dtype=object)
```

Explore total_sqft feature

```
def is_float(x):
    try:
        float(x)
    except:
        return False
    return True
```

```
df3[~df3['total_sqft'].apply(is_float)].head(10)
```

	location	size	total_sqft	bath	price	bhk
30	Yelahanka	4 BHK	2100 - 2850	4.0	186.000	4
122	Hebbal	4 BHK	3067 - 8156	4.0	477.000	4
137	8th Phase JP Nagar	2 BHK	1042 - 1105	2.0	54.005	2
165	Sarjapur	2 BHK	1145 - 1340	2.0	43.490	2
188	KR Puram	2 BHK	1015 - 1540	2.0	56.800	2
410	Kengeri	1 BHK	34.46Sq. Meter	1.0	18.500	1
549	Hennur Road	2 BHK	1195 - 1440	2.0	63.770	2
648	Arekere	9 Bedroom	4125Perch	9.0	265.000	9
661	Yelahanka	2 BHK	1120 - 1145	2.0	48.130	2
672	Bettahalsoor	4 Bedroom	3090 - 5002	4.0	445.000	4

Above shows that total_sqft can be a range (e.g. 2100-2850). For such case we can just take average of min and max value in the range.

There are other cases such as 34.46Sq. Meter which one can convert to square ft using unit conversion. I am going to just drop such corner cases to keep things simple.

```
def convert_sqft_to_num(x):
    tokens = x.split('-')
    if len(tokens) == 2:
        return (float(tokens[0])+float(tokens[1]))/2
    try:
        return float(x)
    except:
        return None
```

```
convert_sqft_to_num('1598')
```

```
1598.0
```

```
convert_sqft_to_num('3067 - 8156')
```

```
5611.5
```

```
convert_sqft_to_num('34.46Sq. Meter')

df4 = df3.copy()
df4['total_sqft'] = df4['total_sqft'].apply(convert_sqft_to_num)
df4.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3
4	Kothanur	2 BHK	1200.0	2.0	51.00	2



For below row, it shows total_sqft as 5611.5 which is an average of the range 3067 - 8156

```
df4.loc[122]

location      Hebbal
size          4 BHK
total_sqft    5611.5
bath          4.0
price         477.0
bhk           4
Name: 122, dtype: object
```

```
(3067 + 8156)/2

5611.5
```

```
df4.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3
4	Kothanur	2 BHK	1200.0	2.0	51.00	2



▼ Feature Engineering

Add new feature called price per square feet

```
df5 = df4.copy()
df5['price_per_sqft'] = df5['price']*100000/df5['total_sqft']
df5.head()
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2	3699.810606
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4	4615.384615
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3	4305.555556
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3	6245.890861
4	Kothanur	2 BHK	1200.0	2.0	51.00	2	4250.000000



```
len(df5['location'].unique())

1304

df5['location'].value_counts()
```

```

Whitefield      534
Sarjapur Road  392
Electronic City 302
Kanakpura Road 266
Thanisandra     233
...
Vidyapeeta      1
Maruthi Extension 1
Okalipura       1
Old Town        1
Abshot Layout   1
Name: location, Length: 1304, dtype: int64

```

Examine locations which is a categorical variable. We need to apply dimensionality reduction technique here to reduce number of locations

```

df5.location = df5.location.apply(lambda x: x.strip())

location_stats = df5.groupby('location')['location'].agg('count').sort_values(ascending=False)
location_stats

location
Whitefield      535
Sarjapur Road   392
Electronic City 304
Kanakpura Road  266
Thanisandra     236
...
1 Giri Nagar    1
Kanakapura Road, 1
Kanakapura main Road 1
Karnataka Shabarimala 1
whitefiled      1
Name: location, Length: 1293, dtype: int64

len(location_stats[location_stats<=10])

1052

```

▼ Dimensionality Reduction

Any location having less than 10 data points should be tagged as "other" location. This way number of categories can be reduced by huge amount. Later on when we do one hot encoding, it will help us with having fewer dummy columns.

```

location_stats_less_than_10 = location_stats[location_stats<=10]
location_stats_less_than_10

location
Basapura      10
1st Block Koramangala 10
Gunjur Palya   10
Kalkere        10
Sector 1 HSR Layout 10
..
1 Giri Nagar   1
Kanakapura Road, 1
Kanakapura main Road 1
Karnataka Shabarimala 1
whitefiled     1
Name: location, Length: 1052, dtype: int64

len(df5.location.unique())

1293

df5.location = df5.location.apply(lambda x: 'other' if x in location_stats_less_than_10 else x)
len(df5.location.unique())

242

df5.head(10)

```

	location	size	total_sqft	bath	price	bhk	price_per_sqft	
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2	3699.810606	
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4	4615.384615	
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3	4305.555556	
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3	6245.890861	
4	Kothanur	2 BHK	1200.0	2.0	51.00	2	4250.000000	
5	Whitefield	2 BHK	1170.0	2.0	38.00	2	3247.863248	
6	Old Airport Road	4 BHK	2732.0	4.0	204.00	4	7467.057101	
7	Rajaji Nagar	4 BHK	3300.0	4.0	600.00	4	18181.818182	
8	Marathahalli	3 BHK	1310.0	3.0	63.25	3	4828.244275	

Outlier Removal Using Business Logic

As a data scientist when you have a conversation with your business manager (who has expertise in real estate), he will tell you that normally square ft per bedroom is 300 (i.e. 2 bhk apartment is minimum 600 sqft. If you have for example 400 sqft apartment with 2 bhk than that seems suspicious and can be removed as an outlier. We will remove such outliers by keeping our minimum threshold per bhk to be 300 sqft

```
df5[df5.total_sqft/df5.bhk<300].head()
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft	
9	other	6 Bedroom	1020.0	6.0	370.0	6	36274.509804	
45	HSR Layout	8 Bedroom	600.0	9.0	200.0	8	33333.333333	
58	Murugeshpalya	6 Bedroom	1407.0	4.0	150.0	6	10660.980810	
68	Devarachikkanahalli	8 Bedroom	1350.0	7.0	85.0	8	6296.296296	
70	other	3 Bedroom	500.0	3.0	100.0	3	20000.000000	

Check above data points. We have 6 bhk apartment with 1020 sqft. Another one is 8 bhk and total sqft is 600. These are clear data errors that can be removed safely

```
df5.shape
```

```
(13246, 7)
```

```
df6 = df5[~(df5.total_sqft/df5.bhk<300)]
```

```
df6.shape
```

```
(12502, 7)
```

Outlier Removal Using Standard Deviation and Mean

```
df6.price_per_sqft.describe()
```

```
count    12456.000000
mean      6308.502826
std       4168.127339
min        267.829813
25%       4210.526316
50%       5294.117647
75%       6916.666667
max      176470.588235
Name: price_per_sqft, dtype: float64
```

Here we find that min price per sqft is 267 rs/sqft whereas max is 12000000, this shows a wide variation in property prices. We should remove outliers per location using mean and one standard deviation

```
def remove_pps_outliers(df):
    df_out = pd.DataFrame()
```

```

for key, subdf in df.groupby('location'):
    m = np.mean(subdf.price_per_sqft)
    st= np.std(subdf.price_per_sqft)
    reduced_df = subdf[(subdf.price_per_sqft>(m-st)) & (subdf.price_per_sqft<=(m+st))]
    df_out = pd.concat([df_out,reduced_df],ignore_index=True)
return df_out

df7 = remove_pps_outliers(df6)
df7.shape

(10241, 7)

```

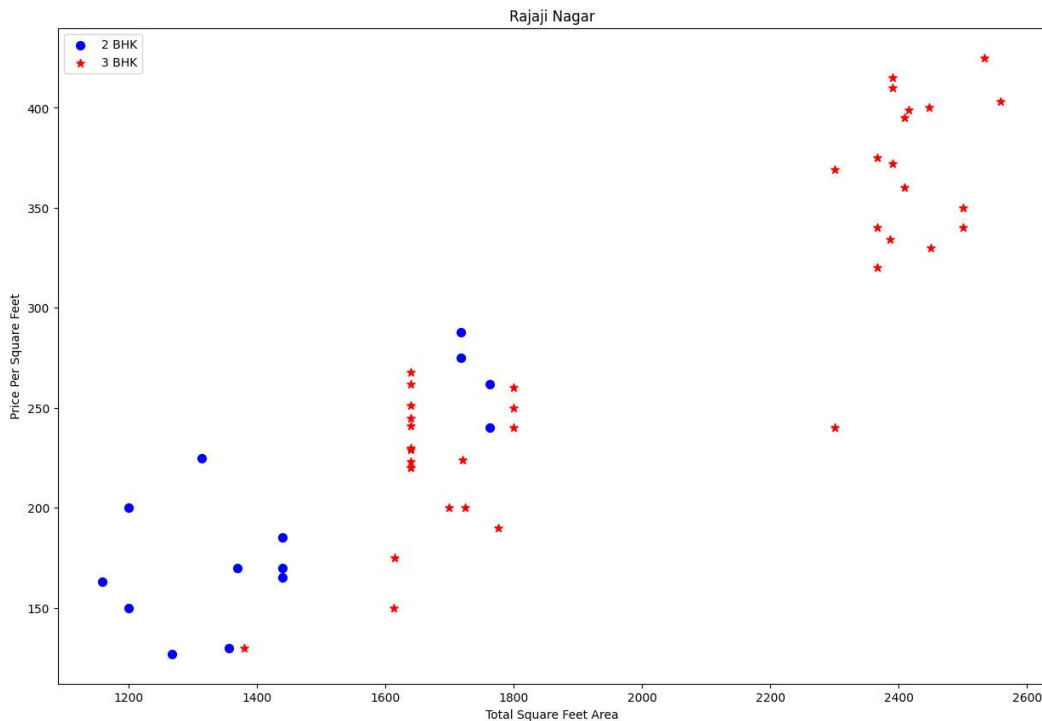
Let's check if for a given location how does the 2 BHK and 3 BHK property prices look like

```

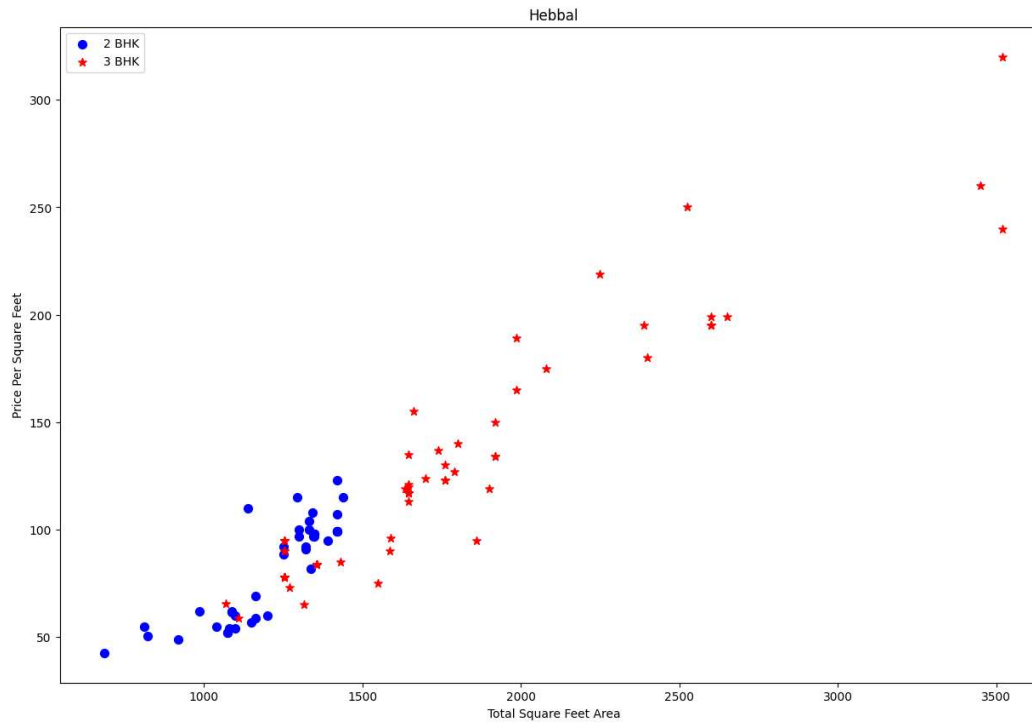
def plot_scatter_chart(df,location):
    bhk2 = df[(df.location==location) & (df.bhk==2)]
    bhk3 = df[(df.location==location) & (df.bhk==3)]
    matplotlib.rcParams['figure.figsize'] = (15,10)
    plt.scatter(bhk2.total_sqft,bhk2.price,color='blue',label='2 BHK', s=50)
    plt.scatter(bhk3.total_sqft,bhk3.price,marker='*', color='red',label='3 BHK', s=50)
    plt.xlabel("Total Square Feet Area")
    plt.ylabel("Price Per Square Feet")
    plt.title(location)
    plt.legend()

plot_scatter_chart(df7,"Rajaji Nagar")

```



```
plot_scatter_chart(df7,"Hebbal")
```



We should also remove properties where for same location, the price of (for example) 3 bedroom apartment is less than 2 bedroom apartment (with same square ft area). What we will do is for a given location, we will build a dictionary of stats per bhk, i.e.

```
{
  '1': {
    'mean': 4000,
    'std': 2000,
    'count': 34
  },
  '2': {
    'mean': 4300,
    'std': 2300,
    'count': 22
  },
}
```

Now we can remove those 2 BHK apartments whose price_per_sqft is less than mean price_per_sqft of 1 BHK apartment.

```
def remove_bhk_outlier(df):
    exclude_indices = np.array([])
    for location, location_df in df.groupby('location'):
        bhk_stats = {}
        for bhk, bhk_df in location_df.groupby('bhk'):
            bhk_stats[bhk] = {
                'mean': np.mean(bhk_df.price_per_sqft),
                'std': np.std(bhk_df.price_per_sqft),
                'count': bhk_df.shape[0]
            }
        # Logic to remove outliers would go here
```



```

for bhk, bhk_df in location_df.groupby('bhk'):
    stats = bhk_stats.get(bhk-1)
    if stats and stats['count']>5:
        exclude_indices = np.append(exclude_indices, bhk_df[bhk_df.price_per_sqft<(stats['mean'])].index.values)
return df.drop(exclude_indices, axis='index')

```

```

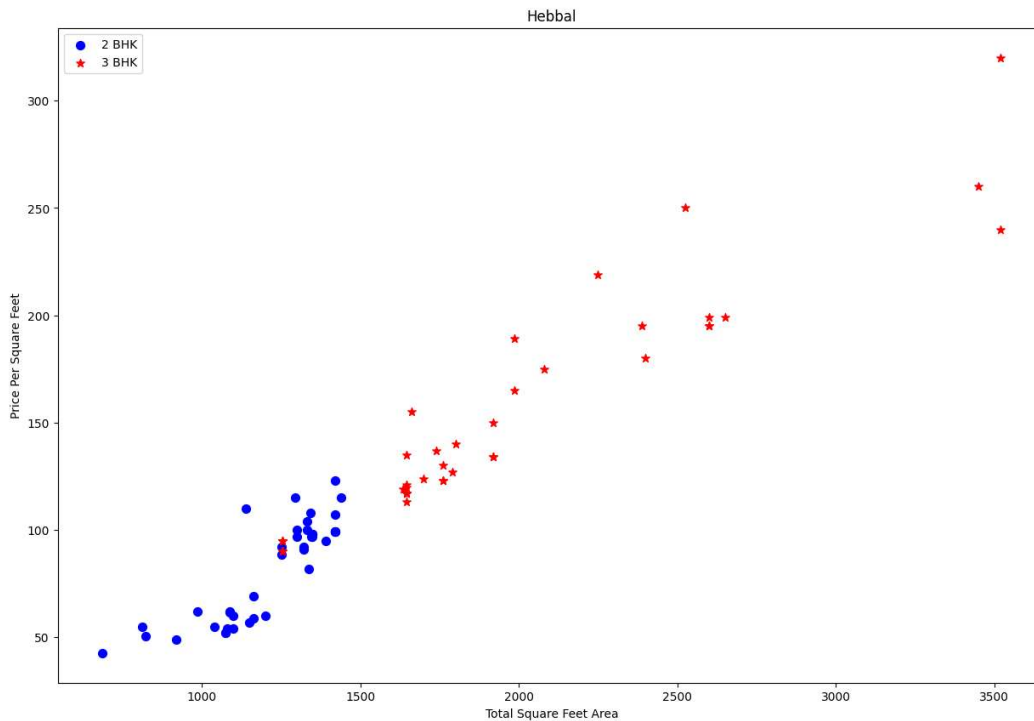
df8 = remove_bhk_outlier(df7)
df8.shape

```

```
(7329, 7)
```

Plot same scatter chart again to visualize price_per_sqft for 2 BHK and 3 BHK properties

```
plot_scatter_chart(df8,"Hebbal")
```



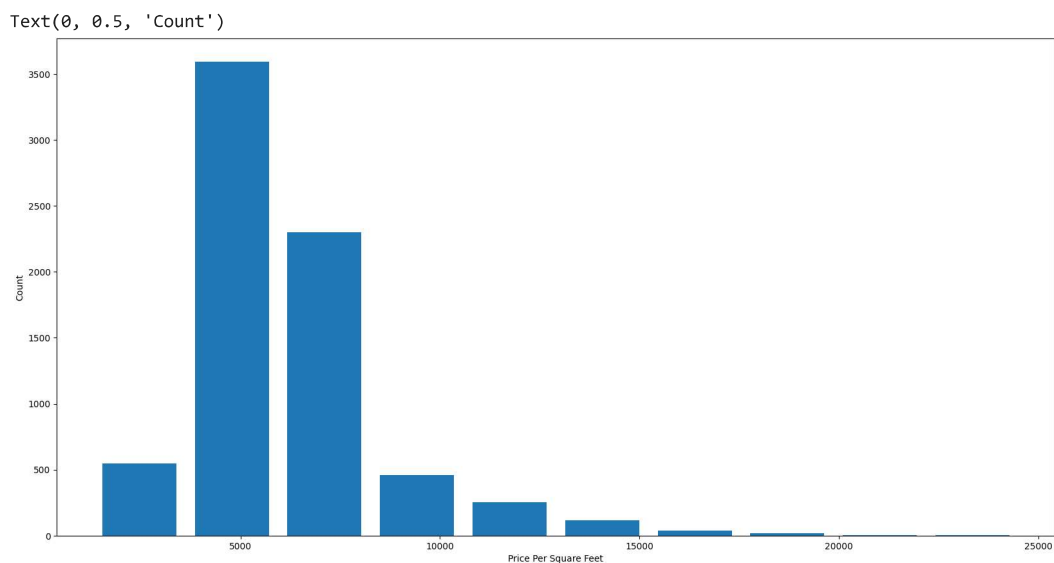
Based on above charts we can see that data points highlighted in red below are outliers and they are being removed due to `remove_bhk_outliers` function

▼ Before and after outlier removal: Hebbal

```

import matplotlib
matplotlib.rcParams["figure.figsize"] = (20,10)
plt.hist(df8.price_per_sqft,rwidth=0.8)
plt.xlabel("Price Per Square Feet")
plt.ylabel("Count")

```



▼ Outlier Removal Using Bathrooms Feature

```
df8.bath.unique()
```

```
array([ 4.,  3.,  2.,  5.,  8.,  1.,  6.,  7.,  9., 12., 16., 13.])
```

```
df8[df8.bath>10]
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft	
5277	Neeladri Nagar	10 BHK	4000.0	12.0	160.0	10	4000.000000	
8486	other	10 BHK	12000.0	12.0	525.0	10	4375.000000	
8575	other	16 BHK	10000.0	16.0	550.0	16	5500.000000	
9308	other	11 BHK	6000.0	12.0	150.0	11	2500.000000	
9639	other	13 BHK	5425.0	13.0	275.0	13	5069.124424	

It is unusual to have 2 more bathrooms than number of bedrooms in a home

```
plt.hist(df8.bath,rwidth=0.8)
plt.xlabel("Number of Bathrooms")
plt.ylabel("Count")
```



```
df8[df8.bath>df8.bhk+2]
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
1626	Chikkabanavar	4 Bedroom	2460.0	7.0	80.0	4	3252.032520
5238	Nagasandra	4 Bedroom	7000.0	8.0	450.0	4	6428.571429
6711	Thanisandra	3 BHK	1806.0	6.0	116.0	3	6423.034330
8411	other	6 BHK	11338.0	9.0	1000.0	6	8819.897689

Again the business manager has a conversation with you (i.e. a data scientist) that if you have 4 bedroom home and even if you have bathroom in all 4 rooms plus one guest bathroom, you will have total bath = total bed + 1 max. Anything above that is an outlier or a data error and can be removed

```
df9 = df8[df8.bath<df8.bhk+2]
```

```
df9.shape
```

(7251, 7)

```
df10 = df9.drop(['size','price_per_sqft'], axis='columns')
df10.head()
```

	location	total_sqft	bath	price	bhk
0	1st Block Jayanagar	2850.0	4.0	428.0	4
1	1st Block Jayanagar	1630.0	3.0	194.0	3
2	1st Block Jayanagar	1875.0	2.0	235.0	3
3	1st Block Jayanagar	1200.0	2.0	130.0	3
4	1st Block Jayanagar	1235.0	2.0	148.0	2

Use One Hot Encoding For Location

```
dummies = pd.get_dummies(df10.location)
dummies.head()
```

	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	6th Phase JP Nagar	7th Phase JP Nagar	8th Phase JP Nagar	9th Phase JP Nagar	...	Vishveshwarya Layout
0	1	0	0	0	0	0	0	0	0	0	...	(
1	1	0	0	0	0	0	0	0	0	0	...	(
2	1	0	0	0	0	0	0	0	0	0	...	(
3	1	0	0	0	0	0	0	0	0	0	...	(
4	1	0	0	0	0	0	0	0	0	0	...	(

5 rows × 242 columns



```
df11 = pd.concat([df10,dummies.drop('other',axis='columns')],axis='columns')
df11.head()
```

	location	total_sqft	bath	price	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	...	Vijayan
0	1st Block Jayanagar	2850.0	4.0	428.0	4	1	0	0	0	0	...	
1	1st Block Jayanagar	1630.0	3.0	194.0	3	1	0	0	0	0	...	
2	1st Block Jayanagar	1875.0	2.0	235.0	3	1	0	0	0	0	...	
3	1st Block Jayanagar	1200.0	2.0	130.0	3	1	0	0	0	0	...	
4	1st Block Jayanagar	1235.0	2.0	148.0	2	1	0	0	0	0	...	

5 rows × 246 columns



```
df12 = df11.drop('location',axis='columns')
df12.head()
```

	total_sqft	bath	price	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	...	Vijayanagar
0	2850.0	4.0	428.0	4	1	0	0	0	0	0	...	(
1	1630.0	3.0	194.0	3	1	0	0	0	0	0	...	(
2	1875.0	2.0	235.0	3	1	0	0	0	0	0	...	(
3	1200.0	2.0	130.0	3	1	0	0	0	0	0	...	(
4	1235.0	2.0	148.0	2	1	0	0	0	0	0	...	(

5 rows × 245 columns



Build a Model Now

```
df12.shape

(7251, 245)

X = df12.drop('price', axis='columns')
X.head()
```

```

total_sqft  bath  bhk  1st Block  1st  2nd  2nd  5th  5th  6th
Jayanagar  Phase  Phase  Stage  Block  Phase  Phase  ...  Vijayanagar
JP  Judicial  Nagarbhavi  Hbr  JP  JP

y = df12.price
y.head()

0    428.0
1    194.0
2    235.0
3    130.0
4    148.0
Name: price, dtype: float64

```

```

# Split the data into training and testing sets
from sklearn.model_selection import train_test_split

```

```

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=12)

```

```

X_train.head()

```

	total_sqft	bath	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	6th Phase JP Nagar	...	Vijayanagar
4425	930.0	1.0	1	0	0	0	0	0	0	0	...	
6126	1418.0	2.0	2	0	0	0	0	0	0	0	...	
9866	1500.0	2.0	3	0	0	0	0	0	0	0	...	
9681	703.0	2.0	2	0	0	0	0	0	0	0	...	
4351	480.0	1.0	1	0	0	0	0	0	0	0	...	

5 rows × 244 columns



```

# Create and train the model with the scaled data
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
lr.fit(X_train, y_train)

```

```

LinearRegression()

```

```

from sklearn.linear_model import LinearRegression

```

```

lr = LinearRegression()
lr.fit(X_train,y_train)
lr.score(X_test,y_test)

```

```

0.8348000128099851

```

```

# Evaluate the model on the training set
train_score = lr.score(X_train, y_train)
print("Training Score:", train_score)

```

```

# Evaluate the model on the testing set
test_score = lr.score(X_test, y_test)
print("Testing Score:", test_score)

```

```

Training Score: 0.8609735915892132
Testing Score: 0.8348000128099851

```

```

from sklearn.metrics import r2_score

```

```

y_pred=lr.predict(X_test)
f'R2_Score:{r2_score(y_test,y_pred)}'

```

```

'R2_Score:0.8348000128099851'

```

```

from sklearn.metrics import mean_squared_error
import math

mse = mean_squared_error(y_test, y_pred)

print("MSE:",mse)

f'RMSE: {math.sqrt(mse)}'

MSE: 1730.762493803334
'RMSE: 41.602433748560124'

```

▼ Use K Fold cross validation to measure accuracy of our LinearRegression model

```

from sklearn.model_selection import ShuffleSplit
from sklearn.model_selection import cross_val_score

cv= ShuffleSplit(n_splits=5, test_size=0.3, random_state=2)

cross_val_score(LinearRegression(), X, y, cv=cv)

array([0.87690805, 0.85547384, 0.8667124 , 0.80282047, 0.86418576])

```

We can see that in 5 iterations we get a score above 80% all the time. This is pretty good but we want to test few other algorithms for regression to see if we can get even better score. We will use GridSearchCV for this purpose

▼ Find best model using GridSearchCV

```

from sklearn.model_selection import GridSearchCV


from sklearn.linear_model import Lasso
from sklearn.tree import DecisionTreeRegressor

def find_best_model_using_gridsearchcv(X,y):
    algos = {
        'linear_regression': {
            'model': LinearRegression(),
            'params': {
            }
        },
        'lasso': {
            'model': Lasso(),
            'params': {
                'alpha': [1,2],
                'selection': ['random', 'cyclic']
            }
        },
        'decision_tree': {
            'model': DecisionTreeRegressor(),
            'params': {
                'criterion': ['mse', 'friedman_mse'],
                'splitter': ['best', 'random']
            }
        }
    }
    scores = []
    cv = ShuffleSplit(n_splits=5, test_size=0.3, random_state=2)
    for algo_name, config in algos.items():
        gs = GridSearchCV(config['model'], config['params'], cv=cv, return_train_score=False)
        gs.fit(X,y)
        scores.append({
            'model': algo_name,
            'best_score': gs.best_score_,
            'best_params': gs.best_params_
        })

    return pd.DataFrame(scores,columns=['model', 'best_score', 'best_params'])

find_best_model_using_gridsearchcv(X,y)

```

	model	best_score	best_params	
0	linear_regression	0.853220	{}	
1	lasso	0.718998	{'alpha': 1, 'selection': 'random'}	
2	decision_tree	0.723622	{'criterion': 'friedman_mse', 'splitter': 'best'}	

Based on above results we can say that LinearRegression gives the best score. Hence we will use that.

▼ Test the model for few properties:

```
def predict_price(location,sqft,bath,bhk):
    loc_index = np.where(X.columns==location)[0][0]

    x = np.zeros(len(X.columns))
    x[0] = sqft
    x[1] = bath
    x[2] = bhk
    if loc_index >= 0:
        x[loc_index] = 1

    return lr.predict([x])[0]
```

```
predict_price('1st Phase JP Nagar',1000, 2, 2)

82.00911568034275
```

```
predict_price('1st Phase JP Nagar',1100, 2, 3)

86.11167949240634
```

```
predict_price('Indira Nagar',1000, 2, 2)

179.2728041910521
```

```
predict_price('Indira Nagar',1200, 3, 3)

193.83271302815285
```

```
predict_price('Indira Nagar',2000, 3, 4)

255.4572453380568
```