

# Data Load: Load banglore home prices into a dataframe

```
In [30]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import matplotlib
matplotlib.rcParams["figure.figsize"] = (20,10)
import warnings
warnings.filterwarnings('ignore')
```

```
In [31]: df1 = pd.read_csv('bengaluru_house_prices.csv')
```

```
In [32]: df1.head()
```

```
Out[32]:
```

	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

```
In [34]: df1.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13320 entries, 0 to 13319
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype  
---  -
0   area_type       13320 non-null  object  
1   availability     13320 non-null  object  
2   location        13319 non-null  object  
3   size            13304 non-null  object  
4   society         7818 non-null   object  
5   total_sqft      13320 non-null  object  
6   bath            13247 non-null  float64 
7   balcony         12711 non-null  float64 
8   price           13320 non-null  float64 
dtypes: float64(3), object(6)
memory usage: 936.7+ KB
```

```
In [33]: df1.shape
```

Out[33]: (13320, 9)

In [35]: `df1.groupby('area_type')['area_type'].agg('count')`

Out[35]:

area_type	
Built-up Area	2418
Carpet Area	87
Plot Area	2025
Super built-up Area	8790

Name: area\_type, dtype: int64

In [36]: `df2 = df1.drop(['area_type', 'society', 'balcony', 'availability'], axis='columns')`  
`df2.head()`

Out[36]:

	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

Drop features that are not required to build our model

## Data Cleaning: Handle NA values

In [37]: `df2.isnull().sum()`

Out[37]:

location	1
size	16
total_sqft	0
bath	73
price	0

dtype: int64

In [38]: `df3=df2.dropna()`  
`df3.isnull().sum()`

Out[38]:

location	0
size	0
total_sqft	0
bath	0
price	0

dtype: int64

In [39]: `df3.shape`

Out[39]: (13246, 5)

## Feature Engineering

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

```
In [12]: df3['size'].unique()
```

```
Out[12]: 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)
```

```
In [40]: df3['bhk'] = df3['size'].apply(lambda x: int(x.split(' ')[0]))
```

```
In [41]: df3.head()
```

```
Out[41]:
```

	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

```
In [42]: df3['bhk'].unique()
```

```
Out[42]: array([ 2,  4,  3,  6,  1,  8,  7,  5, 11,  9, 27, 10, 19, 16, 43, 14, 12,
        13, 18], dtype=int64)
```

## Explore total\_sqft feature

```
In [43]: df3[df3.bhk>20]
```

```
Out[43]:
```

	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

```
In [44]: df3.total_sqft.unique()
```

```
Out[44]: array(['1056', '2600', '1440', ..., '1133 - 1384', '774', '4689'],
        dtype=object)
```

```
In [45]: def is_float(x):
        try:
            float(x)
        except:
            return False
        return True
```

```
In [46]: df3[~df3['total_sqft'].apply(is_float)].head()
```

Out[46]:

	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

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

In [57]:

```
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
```

In [58]:

```
convert_sqft_to_num('2100')
```

Out[58]:

2100.0

In [59]:

```
convert_sqft_to_num('2100 - 2400')
```

Out[59]:

2250.0

In [60]:

```
(2100+2400)/2
```

Out[60]:

2250.0

In [61]:

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

Out[61]:

	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

In [62]:

```
df4.loc[30]
```

```
Out[62]: location      Yelahanka
size              4 BHK
total_sqft       2475.0
bath              4.0
price            186.0
bhk              4
Name: 30, dtype: object
```

## Feature Engineering

### Add new feature called price per square feet

```
In [63]: df5=df4.copy()
df5['price_per_sqft'] = df5['price']*100000/df5['total_sqft']
df5.head()
```

```
Out[63]:
```

	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

```
In [113... df5_stats = df5['price_per_sqft'].describe()
df5_stats
```

```
Out[113]: count      1.320000e+04
mean        7.920759e+03
std         1.067272e+05
min         2.678298e+02
25%         4.267701e+03
50%         5.438331e+03
75%         7.317073e+03
max         1.200000e+07
Name: price_per_sqft, dtype: float64
```

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

```
In [65]: len(df5.location.unique())
```

```
Out[65]: 1304
```

```
In [67]: df5.location = df5.location.apply(lambda x: x.strip())

location_stats = df5.groupby('location')['location'].agg('count').sort_values(ascending=True)
print(location_stats)
```

```

location
Whitefield          535
Sarjapur Road       392
Electronic City      304
Kanakapura 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

```

```
In [68]: len(location_stats[location_stats<=10])
```

```
Out[68]: 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

```
In [69]: location_stats_less_than_10 = location_stats[location_stats<=10]
location_stats_less_than_10
```

```

Out[69]: 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

```

```
In [70]: df5.location = df5.location.apply(lambda x: 'other' if x in location_stats_less_than_10 else len(df5.location.unique()))
```

```
Out[70]: 242
```

```
In [71]: df5.head(10)
```

Out[71]:

	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
9	other	6 Bedroom	1020.0	6.0	370.00	6	36274.509804

## 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

In [72]: `df5[df5.total_sqft/df5.bhk<300].head()`

Out[72]:

	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

In [73]: `df5.shape`

Out[73]: (13246, 7)

In [74]: `df6 = df5[~(df5.total_sqft/df5.bhk<300)]`  
`df6.shape`

Out[74]: (12502, 7)

## Outlier Removal Using Standard Deviation and Mean

```
In [75]: df6.price_per_sqft.describe()
```

```
Out[75]: 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

```
In [76]: 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<=(
    df_out = pd.concat([df_out,reduced_df],ignore_index=True)
return df_out
df7 = remove_pps_outliers(df6)
df7.shape
```

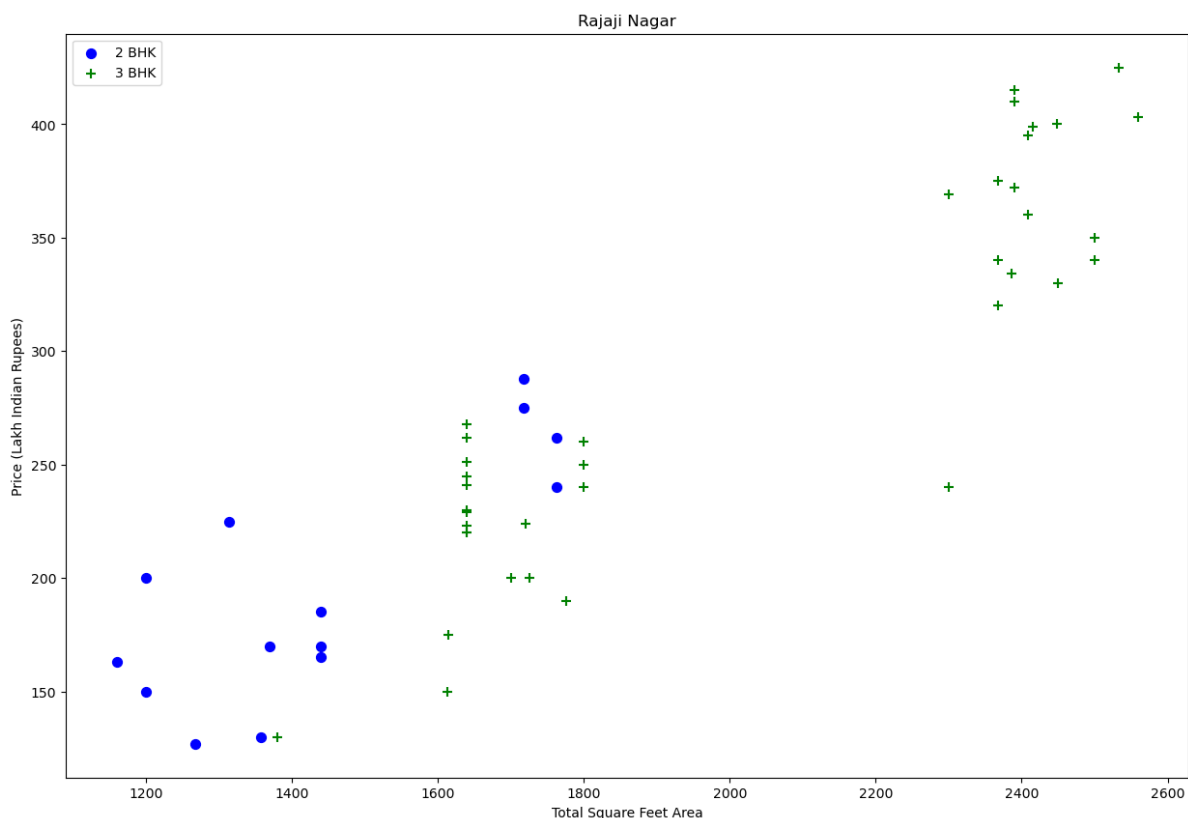
```
Out[76]: (10241, 7)
```

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

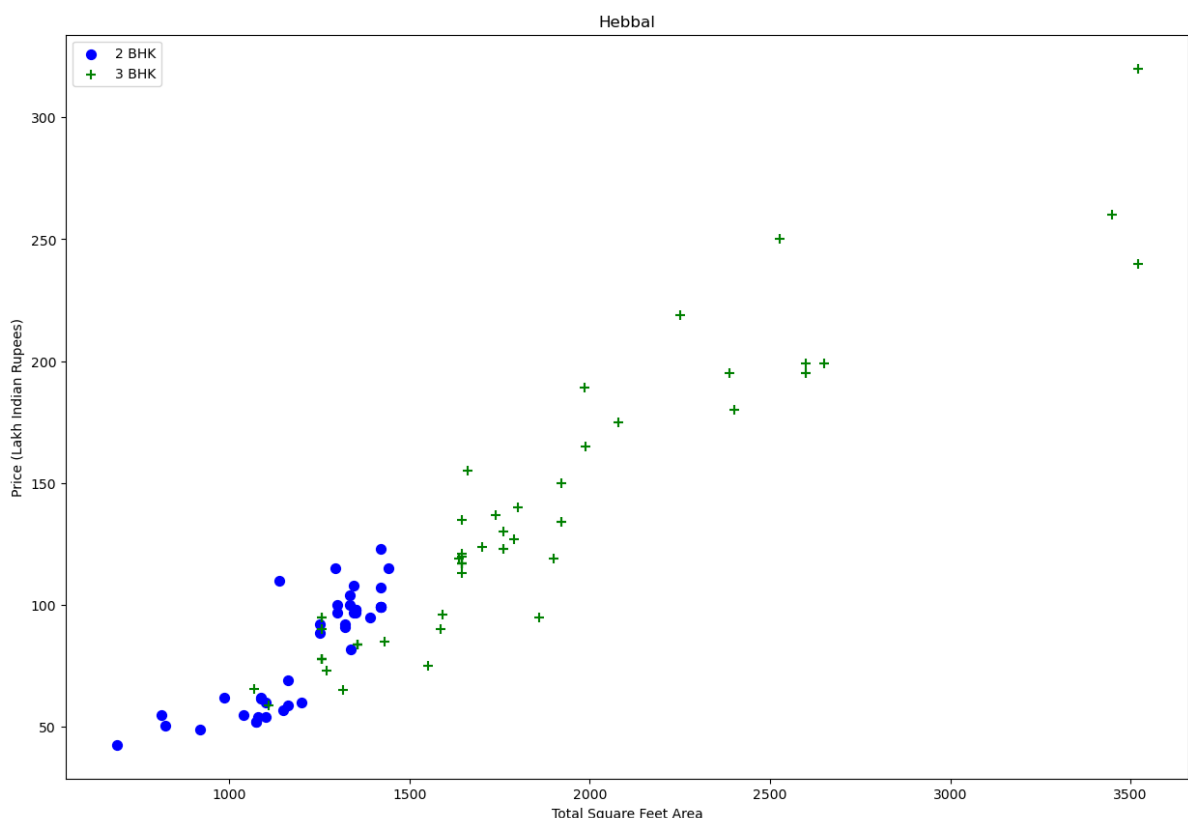
```
In [77]: def plot_scatter_chart(df,location):
bhk2 = df[(df.location==location) & (df.bhk==2)]
bhk3 = df[(df.location==location) & (df.bhk==3)]
matplotlib.rcParams['figure.figsize'] = (15,10)
plt.scatter(bhk2.total_sqft,bhk2.price,color='blue',label='2 BHK', s=50)
plt.scatter(bhk3.total_sqft,bhk3.price,marker='+', color='green',label='3 BHK',
plt.xlabel("Total Square Feet Area")
plt.ylabel("Price (Lakh Indian Rupees)")
plt.title(location)
plt.legend()

plot_scatter_chart(df7,"Rajaji Nagar")
```





```
In [78]: 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 },
}
```

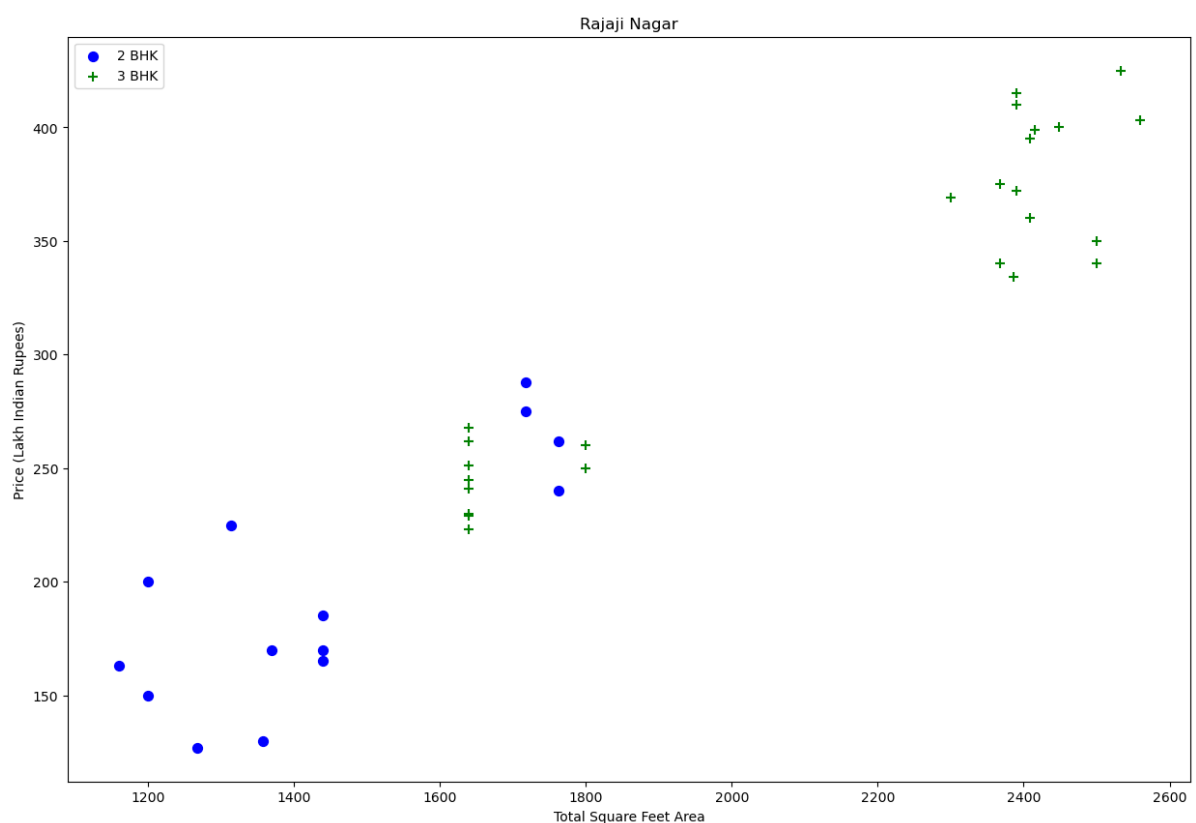
Now we can remove those 2 BHK apartments whose price\_per\_sqft is less than mean price\_per\_sqft of 1 BHK apartment

```
In [79]: def remove_bhk_outliers(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]
        }
    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)
    return df.drop(exclude_indices,axis='index')
df8 = remove_bhk_outliers(df7)
# df8 = df7.copy()
df8.shape
```

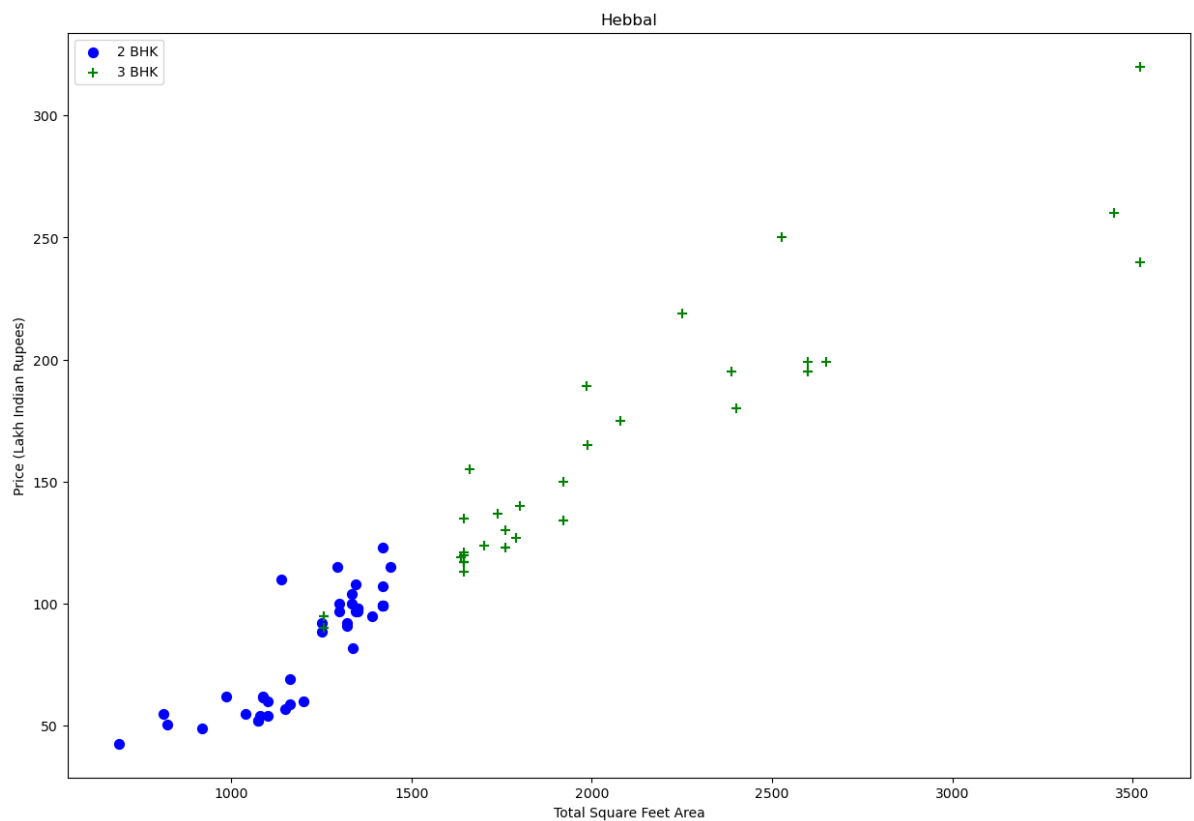
Out[79]: (7329, 7)

Plot same scatter chart again to visualize price\_per\_sqft for 2 BHK and 3 BHK properties

```
In [80]: plot_scatter_chart(df8,"Rajaji Nagar")
```

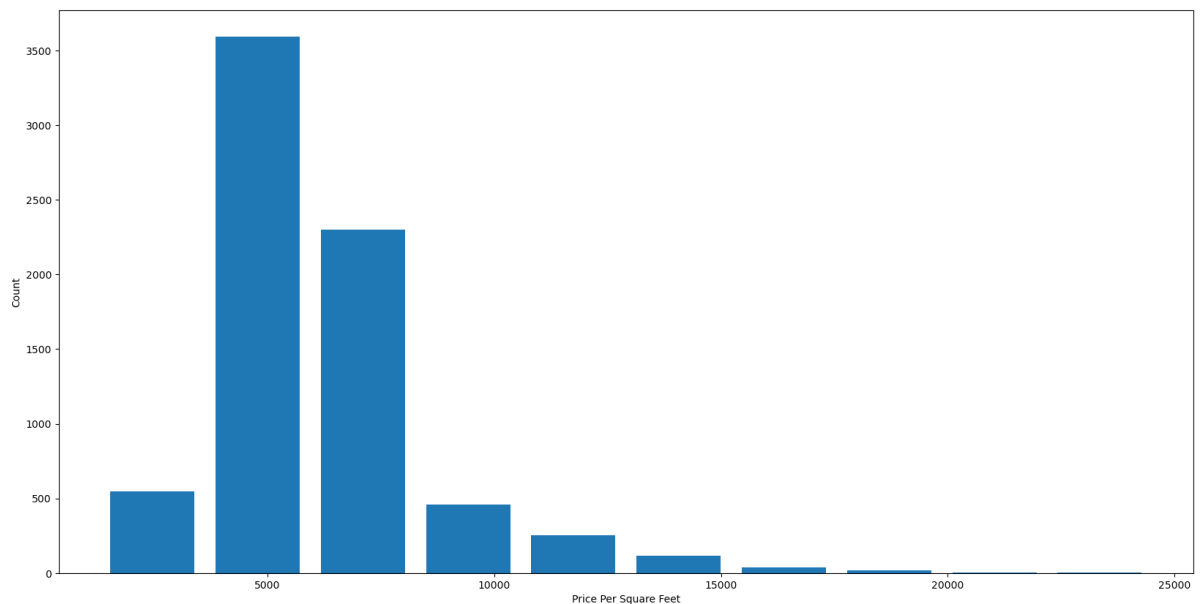


```
In [81]: plot_scatter_chart(df8,"Hebbal")
```



```
In [82]: 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")
```

```
Out[82]: Text(0, 0.5, 'Count')
```



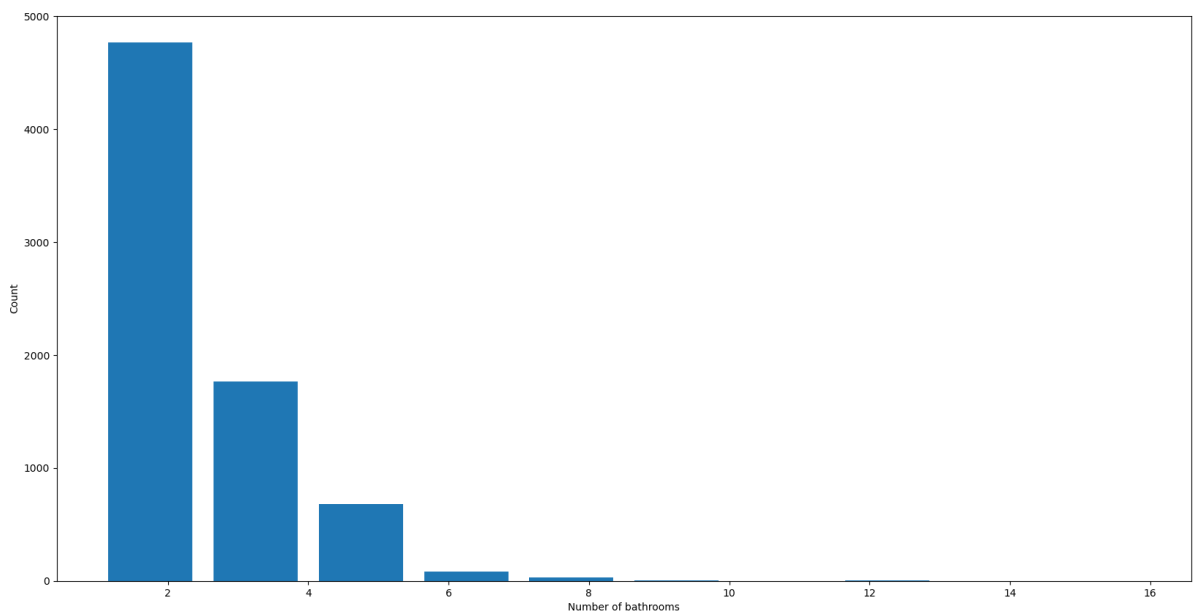
## Outlier Removal Using Bathrooms Feature

```
In [83]: df8.bath.unique()
```

```
Out[83]: array([ 4.,  3.,  2.,  5.,  8.,  1.,  6.,  7.,  9., 12., 16., 13.])
```

```
In [85]: plt.hist(df8.bath,rwidth=0.8)
plt.xlabel("Number of bathrooms")
plt.ylabel("Count")
```

```
Out[85]: Text(0, 0.5, 'Count')
```



```
In [86]: df8[df8.bath>10]
```

```
Out[86]:
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
<b>5277</b>	Neeladri Nagar	10 BHK	4000.0	12.0	160.0	10	4000.000000
<b>8486</b>	other	10 BHK	12000.0	12.0	525.0	10	4375.000000
<b>8575</b>	other	16 BHK	10000.0	16.0	550.0	16	5500.000000
<b>9308</b>	other	11 BHK	6000.0	12.0	150.0	11	2500.000000
<b>9639</b>	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

```
In [87]: df8[df8.bath>df8.bhk+2]
```

```
Out[87]:
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
<b>1626</b>	Chikkabanavar	4 Bedroom	2460.0	7.0	80.0	4	3252.032520
<b>5238</b>	Nagasandra	4 Bedroom	7000.0	8.0	450.0	4	6428.571429
<b>6711</b>	Thanisandra	3 BHK	1806.0	6.0	116.0	3	6423.034330
<b>8411</b>	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

```
In [88]: df9 = df8[df8.bath<df8.bhk+2]
df9.shape
```

```
Out[88]: (7251, 7)
```

```
In [89]: df9.head(2)
```

```
Out[89]:
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
0	1st Block Jayanagar	4 BHK	2850.0	4.0	428.0	4	15017.543860
1	1st Block Jayanagar	3 BHK	1630.0	3.0	194.0	3	11901.840491

```
In [90]: df10 = df9.drop(['size','price_per_sqft'],axis='columns')
df10.head(3)
```

```
Out[90]:
```

	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

## Use One Hot Encoding For Location

```
In [91]: dummies = pd.get_dummies(df10.location)
dummies.head(3)
```

```
Out[91]:
```

	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	...	Vist
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 rows × 242 columns

```
In [92]: df11 = pd.concat([df10,dummies.drop('other',axis='columns')],axis='columns')
df11.head()
```

Out[92]:

	location	total_sqft	bath	price	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	...	Vijaya
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

```
In [93]: df12 = df11.drop('location',axis='columns')
df12.head(2)
```

Out[93]:

	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	...	Vijaya
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 rows × 245 columns

## Build a Model Now...

```
In [94]: df12.shape
```

Out[94]: (7251, 245)

```
In [95]: X = df12.drop(['price'],axis='columns')
X.head(3)
```

Out[95]:

	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	...	Vijaya
0	2850.0	4.0	4	1	0	0	0	0	0	0	...	
1	1630.0	3.0	3	1	0	0	0	0	0	0	...	
2	1875.0	2.0	3	1	0	0	0	0	0	0	...	

3 rows × 244 columns

```
In [96]: X.shape
```

```
Out[96]: (7251, 244)
```

```
In [98]: y=df12.price  
y.head(3)
```

```
Out[98]: 0    428.0  
1    194.0  
2    235.0  
Name: price, dtype: float64
```

```
In [99]: len(y)
```

```
Out[99]: 7251
```

```
In [100... from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=
```

```
In [101... from sklearn.linear_model import LinearRegression  
lr_clf = LinearRegression()  
lr_clf.fit(X_train,y_train)  
lr_clf.score(X_test,y_test)
```

```
Out[101]: 0.8452277697873772
```

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

```
In [102... from sklearn.model_selection import ShuffleSplit  
from sklearn.model_selection import cross_val_score  
  
cv = ShuffleSplit(n_splits=5, test_size=0.2, random_state=0)  
  
cross_val_score(LinearRegression(), X, y, cv=cv)
```

```
Out[102]: array([0.82430186, 0.77166234, 0.85089567, 0.80837764, 0.83653286])
```

We can see that in 4 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

## Test the model for few properties

```
In [106... 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_clf.predict([x])[0]
```

```
In [107... predict_price('1st Phase JP Nagar',1000, 2, 2)
```

```
Out[107]: 83.49904676965245
```

```
In [108... predict_price('1st Phase JP Nagar',1000, 3, 3)
```

```
Out[108]: 86.80519394990591
```

```
In [109... predict_price('Indira Nagar',1000, 2, 2)
```

```
Out[109]: 181.27815484010569
```

```
In [110... predict_price('Indira Nagar',1000, 3, 3)
```

```
Out[110]: 184.5843020203592
```

## Export the tested model to a pickle file

```
In [111... import pickle  
with open('bangalore_home_prices_model.pickle','wb') as f:  
    pickle.dump(lr_clf,f)
```

## Export location and column information to a file that will be useful later on in our prediction application

```
In [112... import json  
columns = {  
    'data_columns' : [col.lower() for col in X.columns]  
}  
with open("columns.json","w") as f:  
    f.write(json.dumps(columns))
```

```
In [ ]:
```