

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 from matplotlib import pyplot as plt
4 %matplotlib inline
5 import matplotlib
6 matplotlib.rcParams["figure.figsize"] = (20,10)
```

Data Load: Load banglore home prices into a dataframe

In [3]:

```
1 df1 = pd.read_csv("Bengaluru_House_Data.csv")
2 df1
```

Out[3]:

	area_type	availability	location	size	society	total_sqft	bath	balcony
0	Super built-up Area	19-Dec	Electronic City Phase II	2 BHK	Coomee	1056	2.0	1.0
1	Plot Area	Ready To Move	Chikka Tirupathi	4 Bedroom	Theanmp	2600	5.0	3.0
2	Built-up Area	Ready To Move	Uttarahalli	3 BHK	NaN	1440	2.0	3.0
3	Super built-up Area	Ready To Move	Lingadheeranahalli	3 BHK	Soiewre	1521	3.0	1.0
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	NaN	1200	2.0	1.0
...
13315	Built-up Area	Ready To Move	Whitefield	5 Bedroom	ArsiaEx	3453	4.0	0.0
13316	Super built-up Area	Ready To Move	Richards Town	4 BHK	NaN	3600	5.0	NaN
13317	Built-up Area	Ready To Move	Raja Rajeshwari Nagar	2 BHK	Mahla T	1141	2.0	1.0
13318	Super built-up Area	18-Jun	Padmanabhanagar	4 BHK	SollyCI	4689	4.0	1.0
13319	Super built-up Area	Ready To Move	Doddathoguru	1 BHK	NaN	550	1.0	1.0

13320 rows × 9 columns



In [4]:

```
1 df1['area_type'].unique()
```

Out[4]:

```
array(['Super built-up Area', 'Plot Area', 'Built-up Area',  
      'Carpet Area'], dtype=object)
```

In [5]:

```
1 df1['area_type'].value_counts()
```

Out[5]:

```
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

In [6]:

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

Out[6]:

	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
...
13315	Whitefield	5 Bedroom	3453	4.0	231.00
13316	Richards Town	4 BHK	3600	5.0	400.00
13317	Raja Rajeshwari Nagar	2 BHK	1141	2.0	60.00
13318	Padmanabhanagar	4 BHK	4689	4.0	488.00
13319	Doddathoguru	1 BHK	550	1.0	17.00

13320 rows × 5 columns

Data Cleaning: Handle NA values

In [7]:

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

Out[7]:

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

In [8]:

```
1 df3 = df2.dropna()
2 df3.isnull().sum()
```

Out[8]:

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

In [9]:

```
1 df3.shape
```

Out[9]:

```
(13246, 5)
```

Feature Engineering

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

In [10]:

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

<ipython-input-10-66d19fd42217>:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

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

Out[10]:

	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
...
13315	Whitefield	5 Bedroom	3453	4.0	231.00	5
13316	Richards Town	4 BHK	3600	5.0	400.00	4
13317	Raja Rajeshwari Nagar	2 BHK	1141	2.0	60.00	2
13318	Padmanabhanagar	4 BHK	4689	4.0	488.00	4
13319	Doddathoguru	1 BHK	550	1.0	17.00	1

13246 rows × 6 columns

Explore total_sqft feature

In [11]:

```
1 def is_float(x):
2     try:
3         float(x)
4     except:
5         return False
6     return True
```

In [12]:

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

Out[12]:

	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

In [13]:

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

In [14]:

```
1 convert_sqft_to_num('1120 - 1145')
```

Out[14]:

1132.5

In [15]:

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

In [16]:

```

1 df4 = df3.copy()
2 df4['total_sqft'] = df4['total_sqft'].apply(convert_sqft_to_num)
3 df4

```

Out[16]:

	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
...
13315	Whitefield	5 Bedroom	3453.0	4.0	231.00	5
13316	Richards Town	4 BHK	3600.0	5.0	400.00	4
13317	Raja Rajeshwari Nagar	2 BHK	1141.0	2.0	60.00	2
13318	Padmanabhanagar	4 BHK	4689.0	4.0	488.00	4
13319	Doddathoguru	1 BHK	550.0	1.0	17.00	1

13246 rows × 6 columns

Feature Engineering

Add new feature called price per square feet

In [17]:

```

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

```

Out[17]:

	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 [18]:

```
1 df5['price_per_sqft'].describe()
```

Out[18]:

```
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 [19]:

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

In [20]:

```
1 location_stats = df5['location'].value_counts(ascending=False)
2 location_stats
```

Out[20]:

```
Whitefield          535
Sarjapur Road       392
Electronic City     304
Kanakpura Road      266
Thanisandra         236
...
Iggalur             1
12th cross srinivas nagar banshankari 3rd stage  1
P&T Colony          1
Annapoorneshwari Layout, JP nagar 7th phase     1
Beml layout, Rajarajeshwari nagar                1
Name: location, Length: 1293, dtype: int64
```

In [21]:

```
1 location_stats = df5.groupby('location')['location'].count().sort_values(ascending=False)
2 location_stats
```

Out[21]:

```
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
```

In [22]:

```
1 location_stats.values.sum()
```

Out[22]:

13246

In [23]:

```
1 len(location_stats[location_stats>10])
```

Out[23]:

241

In [24]:

```
1 len(location_stats)
```

Out[24]:

1293

In [25]:

```
1 len(location_stats[location_stats<=10])
```

Out[25]:

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 [26]:

```
1 location_stats_less_than_10 = location_stats[location_stats<=10]
2 location_stats_less_than_10
```

Out[26]:

```
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 [27]:

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

Out[27]:

242

In [28]:

```
1 df5.head(10)
```

Out[28]:

	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

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 [29]:

```
1 df5[df5.total_sqft/df5.bhk<300]
```

Out[29]:

	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
...
13277	other	7 Bedroom	1400.0	7.0	218.0	7	15571.428571
13279	other	6 Bedroom	1200.0	5.0	130.0	6	10833.333333
13281	Margondanahalli	5 Bedroom	1375.0	5.0	125.0	5	9090.909091
13303	Vidyaranyapura	5 Bedroom	774.0	5.0	70.0	5	9043.927649
13311	Ramamurthy Nagar	7 Bedroom	1500.0	9.0	250.0	7	16666.666667

744 rows × 7 columns

In [30]:

```
1 df5.shape
```

Out[30]:

(13246, 7)

In [31]:

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

Out[31]:

(12502, 7)

In [32]:

```
1 df6.head(10)
```

Out[32]:

	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
10	Whitefield	3 BHK	1800.0	2.0	70.00	3	3888.888889

Outlier Removal Using Standard Deviation and Mean

In [33]:

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

Out[33]:

```
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 [34]:

```

1 def remove_pps_outliers(df):
2     df_out = pd.DataFrame()
3     for key, subdf in df.groupby('location'):
4         m = np.mean(subdf.price_per_sqft)
5         st = np.std(subdf.price_per_sqft)
6         reduced_df = subdf[(subdf.price_per_sqft>(m-st)) & (subdf.price_per_sqft<=(m+st))
7         df_out = pd.concat([df_out,reduced_df],ignore_index=True)
8     return df_out
9 df7 = remove_pps_outliers(df6)
10 df7.shape

```

Out[34]:

(10241, 7)

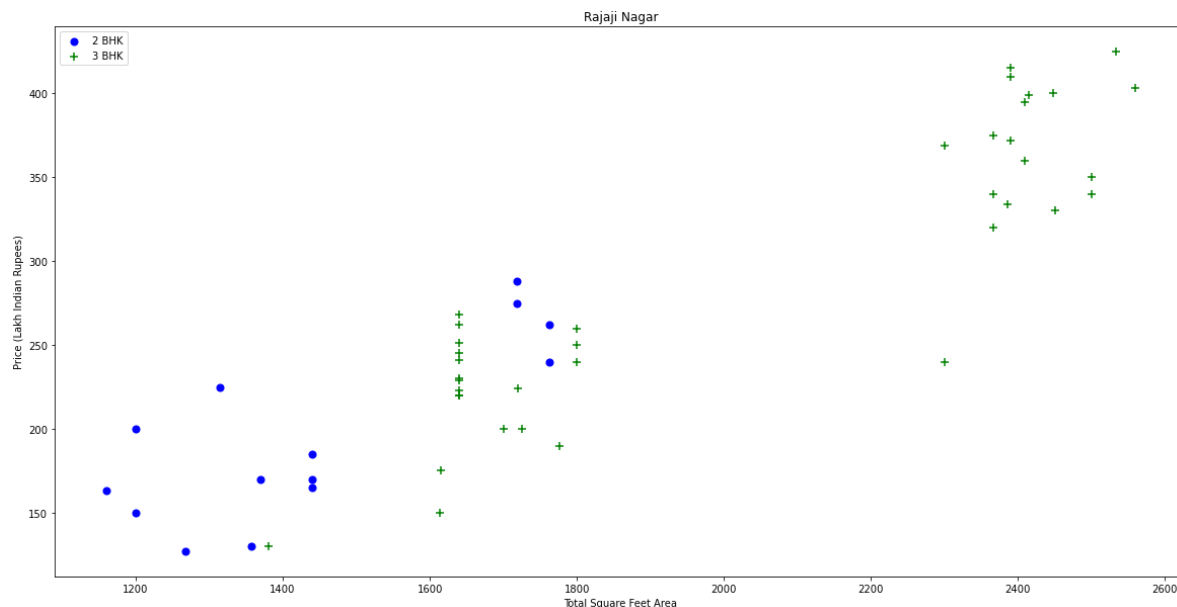
Let's check if for a given location how does the 2 BHK and 3 BHK property prices look like

In [35]:

```

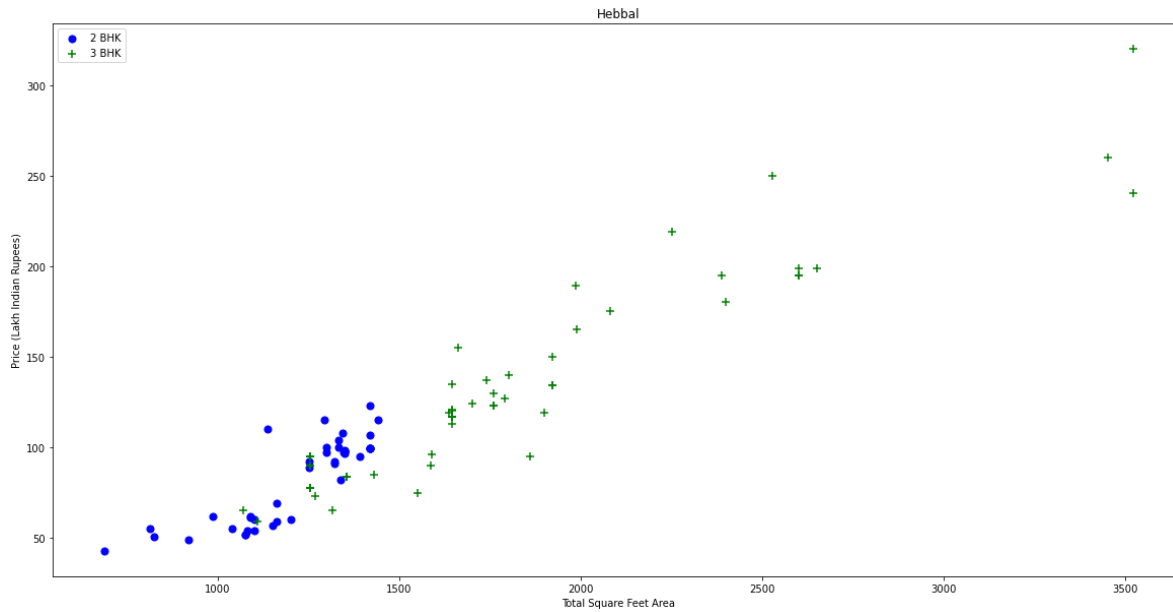
1 def plot_scatter_chart(df,location):
2     bhk2 = df[(df.location==location) & (df.bhk==2)]
3     bhk3 = df[(df.location==location) & (df.bhk==3)]
4     plt.scatter(bhk2.total_sqft,bhk2.price,color='blue',label='2 BHK', s=50)
5     plt.scatter(bhk3.total_sqft,bhk3.price,marker='+', color='green',label='3 BHK', s=50)
6     plt.xlabel("Total Square Feet Area")
7     plt.ylabel("Price (Lakh Indian Rupees)")
8     plt.title(location)
9     plt.legend()
10
11 plot_scatter_chart(df7,"Rajaji Nagar")

```



In [36]:

```
1 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 [37]:

```

1 def remove_bhk_outliers(df):
2     exclude_indices = np.array([])
3     for location, location_df in df.groupby('location'):
4         bhk_stats = {}
5         for bhk, bhk_df in location_df.groupby('bhk'):
6             bhk_stats[bhk] = {
7                 'mean': np.mean(bhk_df.price_per_sqft),
8                 'std': np.std(bhk_df.price_per_sqft),
9                 'count': bhk_df.shape[0]
10            }
11         for bhk, bhk_df in location_df.groupby('bhk'):
12             stats = bhk_stats.get(bhk-1)
13             if stats and stats['count'] > 5:
14                 exclude_indices = np.append(exclude_indices, bhk_df[bhk_df.price_per_sqft > stats['std'] * 3].index)
15     return df.drop(exclude_indices, axis='index')
16 df8 = remove_bhk_outliers(df7)
17 # df8 = df7.copy()
18 df8.shape

```

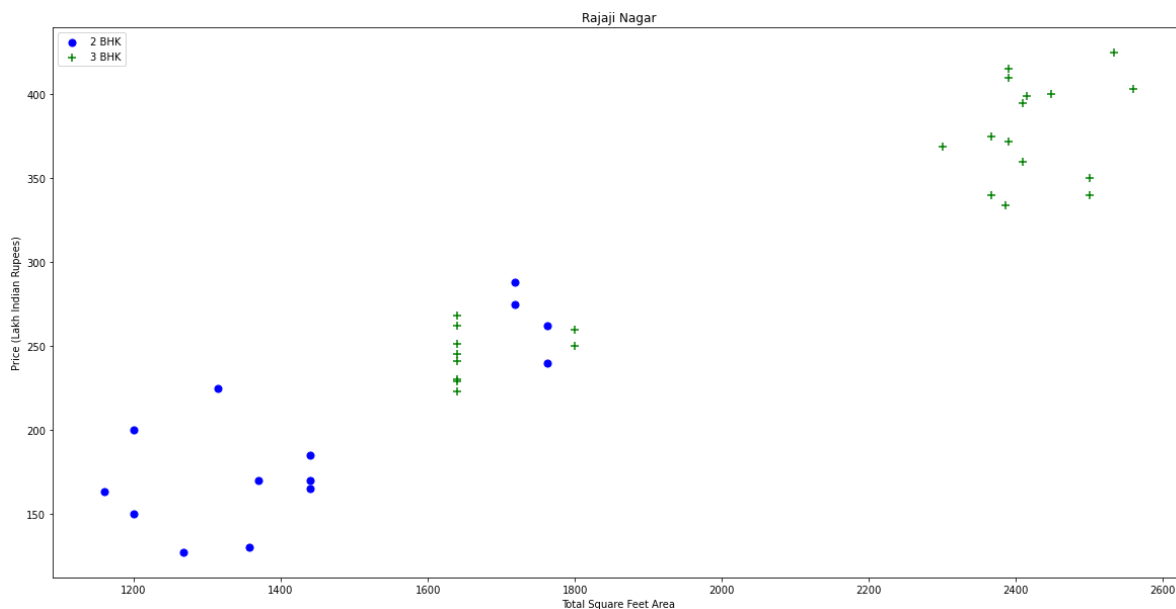
Out[37]:

(7329, 7)

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

In [38]:

```
1 plot_scatter_chart(df8, "Rajaji Nagar")
```

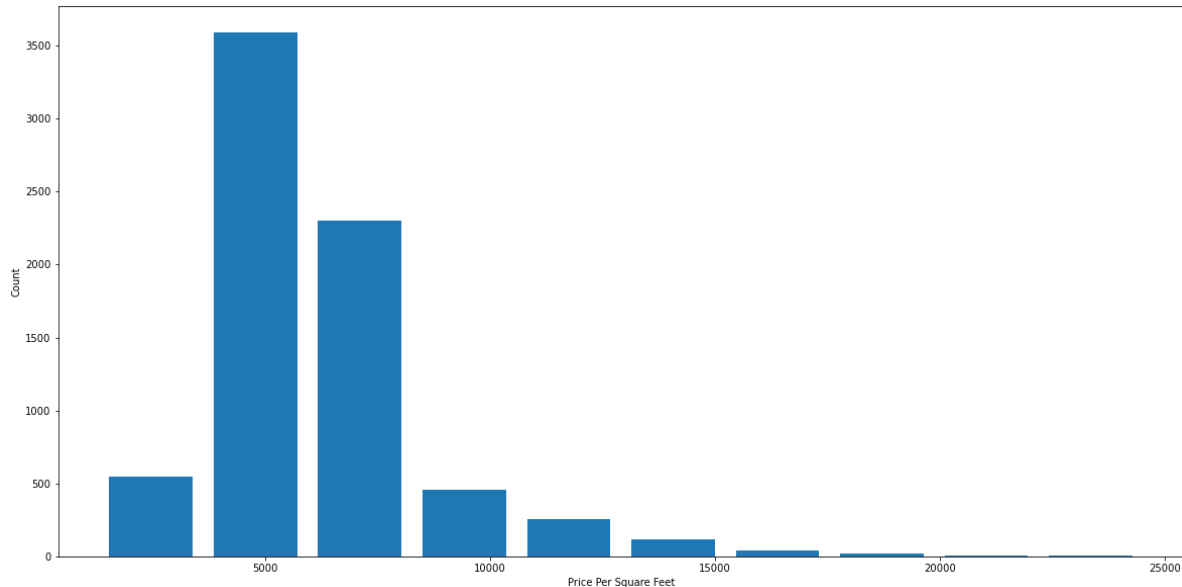


In [39]:

```
1 import matplotlib.pyplot as plt
2 plt.hist(df8.price_per_sqft,rwidth=0.8)
3 plt.xlabel("Price Per Square Feet")
4 plt.ylabel("Count")
```

Out[39]:

Text(0, 0.5, 'Count')



Outlier Removal Using Bathrooms Feature

In [40]:

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

Out[40]:

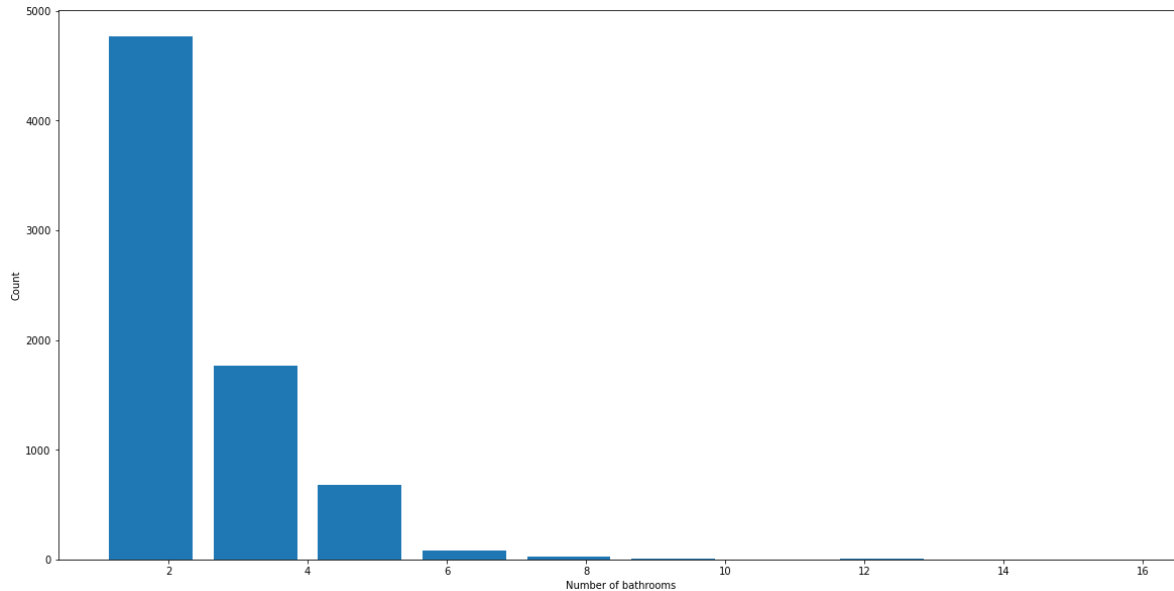
	location	size	total_sqft	bath	price	bhk	price_per_sqft
5277	Neeladri Nagar	10 BHK	4000.0	12.0	160.0	10	4000.000000
8486	other	10 BHK	12000.0	12.0	525.0	10	4375.000000
8575	other	16 BHK	10000.0	16.0	550.0	16	5500.000000
9308	other	11 BHK	6000.0	12.0	150.0	11	2500.000000
9639	other	13 BHK	5425.0	13.0	275.0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

In [41]:

```
1 plt.hist(df8.bath,rwidth=0.8)
2 plt.xlabel("Number of bathrooms")
3 plt.ylabel("Count")
```

Out[41]:

Text(0, 0.5, 'Count')



In [42]:

```
1 df8.shape
```

Out[42]:

(7329, 7)

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

In [43]:

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

Out[43]:

(4, 7)

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 [44]:

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

Out[44]:

(7251, 7)

In [45]:

```
1 df9.head(2)
```

Out[45]:

	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 [46]:

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

Out[46]:

	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
...
10232	other	1200.0	2.0	70.0	2
10233	other	1800.0	1.0	200.0	1
10236	other	1353.0	2.0	110.0	2
10237	other	812.0	1.0	26.0	1
10240	other	3600.0	5.0	400.0	4

7251 rows × 5 columns

Use One Hot Encoding For Location

In [47]:

```
1 dummies = pd.get_dummies(df10.location)
2 dummies
```

Out[47]:

	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
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
...
10232	0	0	0	0	0	0	0	0	0	0
10233	0	0	0	0	0	0	0	0	0	0
10236	0	0	0	0	0	0	0	0	0	0
10237	0	0	0	0	0	0	0	0	0	0
10240	0	0	0	0	0	0	0	0	0	0

7251 rows × 242 columns

In [48]:

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

Out[48]:

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

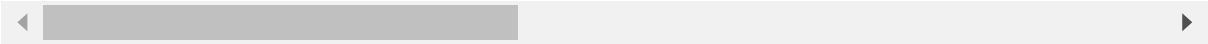
In [49]:

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

Out[49]:

	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	...
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	...
...
10232	1200.0	2.0	70.0	2	0	0	0	0	0	0	...
10233	1800.0	1.0	200.0	1	0	0	0	0	0	0	...
10236	1353.0	2.0	110.0	2	0	0	0	0	0	0	...
10237	812.0	1.0	26.0	1	0	0	0	0	0	0	...
10240	3600.0	5.0	400.0	4	0	0	0	0	0	0	...

7251 rows × 245 columns



Build a Model Now...

In [50]:

```
1 df12.shape
```

Out[50]:

(7251, 245)

In [51]:

```
1 X = df12.drop(['price'],axis='columns')
2 X
```

Out[51]:

	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	.
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	1200.0	2.0	3	1	0	0	0	0	0	0	.
4	1235.0	2.0	2	1	0	0	0	0	0	0	.
...
10232	1200.0	2.0	2	0	0	0	0	0	0	0	.
10233	1800.0	1.0	1	0	0	0	0	0	0	0	.
10236	1353.0	2.0	2	0	0	0	0	0	0	0	.
10237	812.0	1.0	1	0	0	0	0	0	0	0	.
10240	3600.0	5.0	4	0	0	0	0	0	0	0	.

7251 rows × 244 columns

In [52]:

```
1 y = df12.price
2 y.head(3)
```

Out[52]:

```
0    428.0
1    194.0
2    235.0
Name: price, dtype: float64
```

In [53]:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=10)
```

In [54]:

```
1 from sklearn.linear_model import LinearRegression
2 lr_clf = LinearRegression()
3 lr_clf.fit(X_train,y_train)
4 lr_clf.score(X_test,y_test)
```

Out[54]:

0.845227769787429

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

In [60]:

```
1 from sklearn.model_selection import ShuffleSplit
2 from sklearn.model_selection import cross_val_score
3
4 cv = ShuffleSplit(n_splits=5, test_size=0.2, random_state=0)
5
6 cross_val_score(LinearRegression(), X, y, cv=cv)
```

Out[60]:

array([0.82430186, 0.77166234, 0.85089567, 0.80837764, 0.83653286])

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

Test the model for few properties

In [63]:

```
1 def predict_price(location,sqft,bath,bhk):
2     loc_index = np.where(X.columns==location)[0][0]
3
4     x = np.zeros(len(X.columns))
5     x[0] = sqft
6     x[1] = bath
7     x[2] = bhk
8     if loc_index >= 0:
9         x[loc_index] = 1
10
11     return lr_clf.predict([x])[0]
```

In [64]:

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

Out[64]:

83.49904677172407

In [65]:

```
1 predict_price('1st Phase JP Nagar',1000, 3, 3)
```

Out[65]:

86.8051939519899

In [66]:

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

Out[66]:

181.27815484006965

In [67]:

```
1 predict_price('Indira Nagar',1000, 3, 3)
```

Out[67]:

184.58430202033554

Export the tested model to a pickle file

In [68]:

```
1 import pickle
2 with open('bangalore_home_prices_model.pickle','wb') as f:
3     pickle.dump(lr_clf,f)
```

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

In [71]:

```
1 import json
2 columns = {
3     'data_columns' : [col.lower() for col in X.columns]
4 }
5 with open("columns.json",'w') as f:
6     f.write(json.dumps(columns))
```

In []:

```
1
```

