

```

➡ Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client\_id=947:
Enter your authorization code:
.....
Mounted at /content/drive

```

```
df1 = pd.read_csv("/content/drive/My Drive/datasets_20710_26737_Bengaluru_House_Data.csv")
df1.head()
```

| | area_type | availability | location | size | society | total_sqft | bat |
|---|---------------------|---------------|--------------------------|-----------|---------|------------|-----|
| 0 | Super built-up Area | 19-Dec | Electronic City Phase II | 2 BHK | Coomee | 1056 | 2.0 |
| 1 | Plot Area | Ready To Move | Chikka Tirupathi | 4 Bedroom | Theanmp | 2600 | 5.0 |
| 2 | Built-up Area | Ready To Move | Uttarahalli | 3 BHK | NaN | 1440 | 2.0 |
| 3 | Super built-up Area | Ready To Move | Chikka Tirupathi | 4 Bedroom | Theanmp | 2600 | 5.0 |



```
<bound method DataFrame.info of
0      Super built-up Area      19-Dec ...      1.0      39.07
1              Plot Area Ready To Move ...      3.0     120.00
2      Built-up Area Ready To Move ...      3.0      62.00
3      Super built-up Area Ready To Move ...      1.0      95.00
4      Super built-up Area Ready To Move ...      1.0      51.00
...
13315      Built-up Area Ready To Move ...      0.0     231.00
13316      Super built-up Area Ready To Move ...      NaN     400.00
13317      Built-up Area Ready To Move ...      1.0      60.00
13318      Super built-up Area Ready To Move ...      1.0     100.00
```

```
df1.describe()
```

```
↳
```

| | bath | balcony | price |
|--------------|--------------|--------------|--------------|
| count | 13247.000000 | 12711.000000 | 13320.000000 |
| mean | 2.692610 | 1.584376 | 112.565627 |
| std | 1.341458 | 0.817263 | 148.971674 |
| min | 1.000000 | 0.000000 | 8.000000 |
| 25% | 2.000000 | 1.000000 | 50.000000 |
| 50% | 2.000000 | 2.000000 | 72.000000 |
| 75% | 3.000000 | 2.000000 | 120.000000 |
| max | 40.000000 | 3.000000 | 3600.000000 |

```
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
```

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

```
↳ (13320, 5)
```

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

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

```
df2.shape
```

```
↳ (13320, 5)
```

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

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

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

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py: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/10min.html>
 """Entry point for launching an IPython kernel.

```
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
```

```
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)
```

```
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 |

```

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

```

```

df4 = df3.copy()
df4.total_sqft = df4.total_sqft.apply(convert_sqft_to_num)
df4 = df4[df4.total_sqft.notnull()]
df4.head(2)

```

| | 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 |

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

```

location    Yelahanka
size        4 BHK
total_sqft  2475
bath        4
price       186
bhk         4
Name: 30, dtype: object

```

```
(2100+2850)/2
```

```
↳ 2475.0
```

```
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 |

```
df5.to_csv("bhp.csv",index=False)
```

```
df5.location = df5.location.apply(lambda x: x.strip())
location_stats = df5['location'].value_counts(ascending=False)
location_stats
```

```
↳
```

| | |
|------------------------------|-----|
| Whitefield | 533 |
| Sarjapur Road | 392 |
| Electronic City | 304 |
| Kanakpura Road | 264 |
| Thanisandra | 235 |
| ... | |
| Binny Mills Employees Colony | 1 |
| Yemlur, Old Airport Road, | 1 |
| N R Layout | 1 |
| Tharabanahalli | 1 |
| Uvce Layout | 1 |

Name: location, Length: 1287, dtype: int64

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

```
↳ 13200
```

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

```
↳ 240
```

```
len(location_stats)
```

```
↳ 1287
```

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

```
↳ 1047
```

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

```
↳
BTM 1st Stage          10
1st Block Koramangala  10
Naganathapura          10
Nagappa Reddy Layout   10
Dodsworth Layout      10
..
Binny Mills Employees Colony  1
Yemlur, Old Airport Road,    1
N R Layout              1
Tharabanahalli          1
Uvce Layout             1
Name: location, Length: 1047, dtype: int64
```

```
len(df5.location.unique())
```

```
↳ 1287
```

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

```
↳ 241
```

```
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 |
| 9 | other | 6 Bedroom | 1020.0 | 6.0 | 370.00 | 6 | 36274.509804 |

```
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 |
| 50 | Munambakkam | 6 Bedroom | 1107.0 | 4.0 | 150.0 | 6 | 13558.174345 |

df5.shape

```
(13200, 7)
```

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

df6.shape

```
(12456, 7)
```

```
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
```

```
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

```
(10242, 7)
```

```
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', s=50)
```

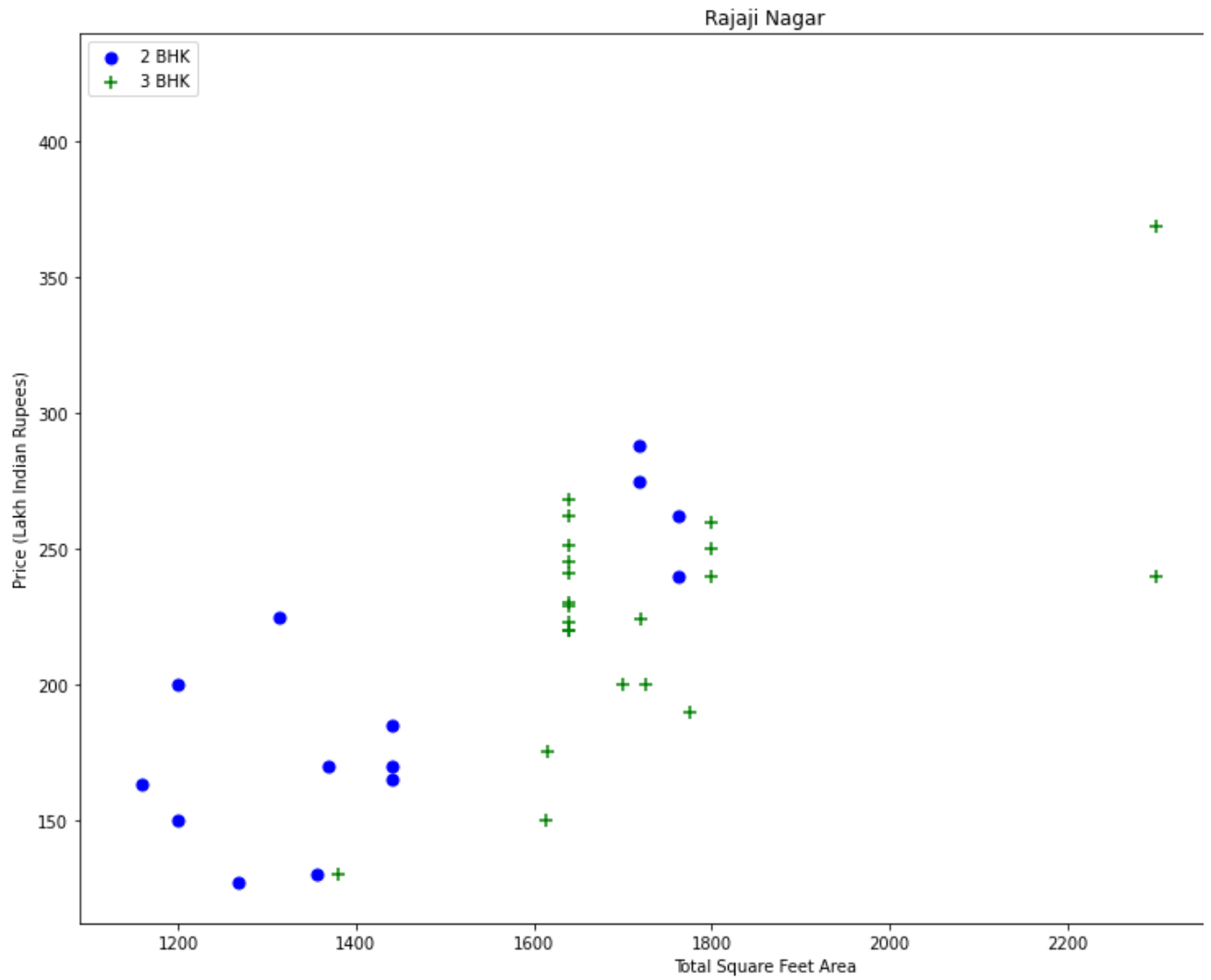
```
    plt.xlabel("Total Square Feet Area")
```

```
    plt.ylabel("Price (Lakh Indian Rupees)")
```

```
    plt.title(location)
```

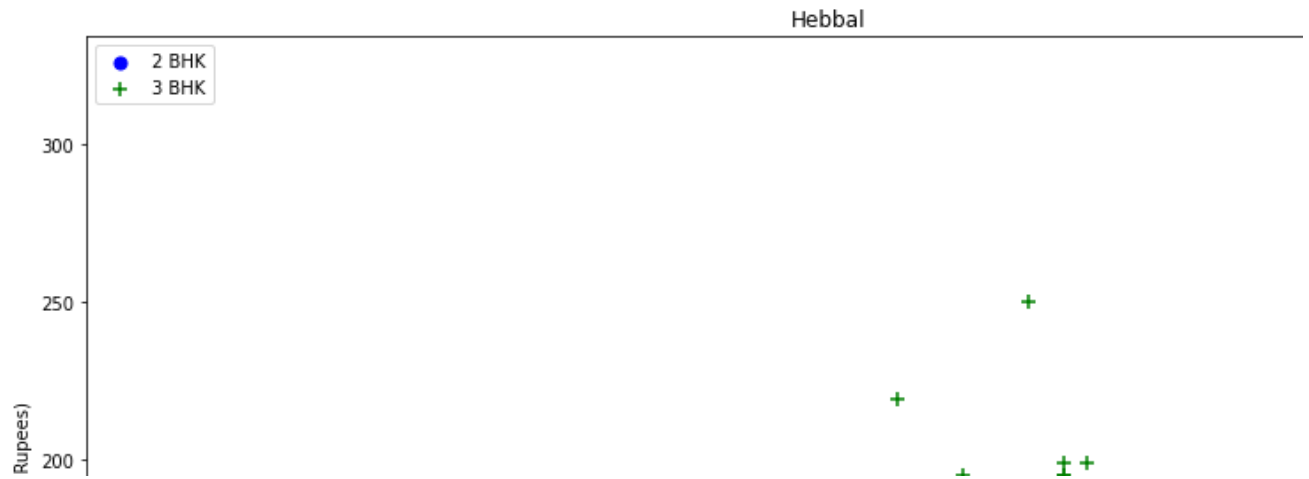
```
    plt.legend()
```

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



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



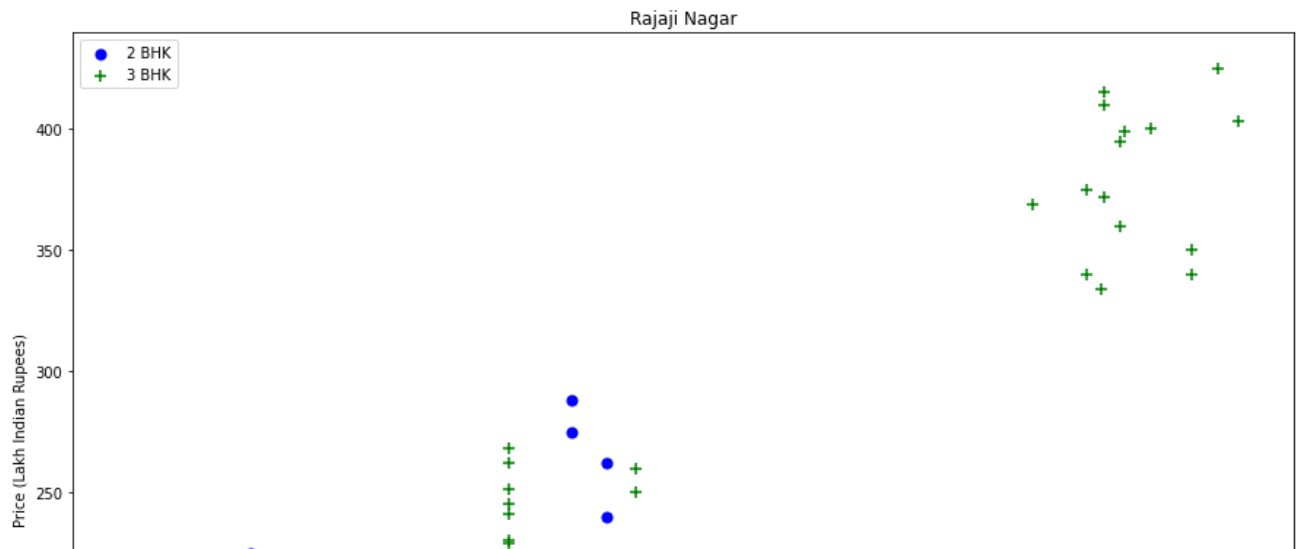


```
def remove_bhk_outliers(df):
    exclude_indices = np.array([])
    for location, location_df in df.groupby('location'):
        bhk_stats = {}
        for bhk, bhk_df in location_df.groupby('bhk'):
            bhk_stats[bhk] = {
                'mean': np.mean(bhk_df.price_per_sqft),
                '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<
            return df.drop(exclude_indices,axis='index')
df8 = remove_bhk_outliers(df7)
# df8 = df7.copy()
df8.shape
```

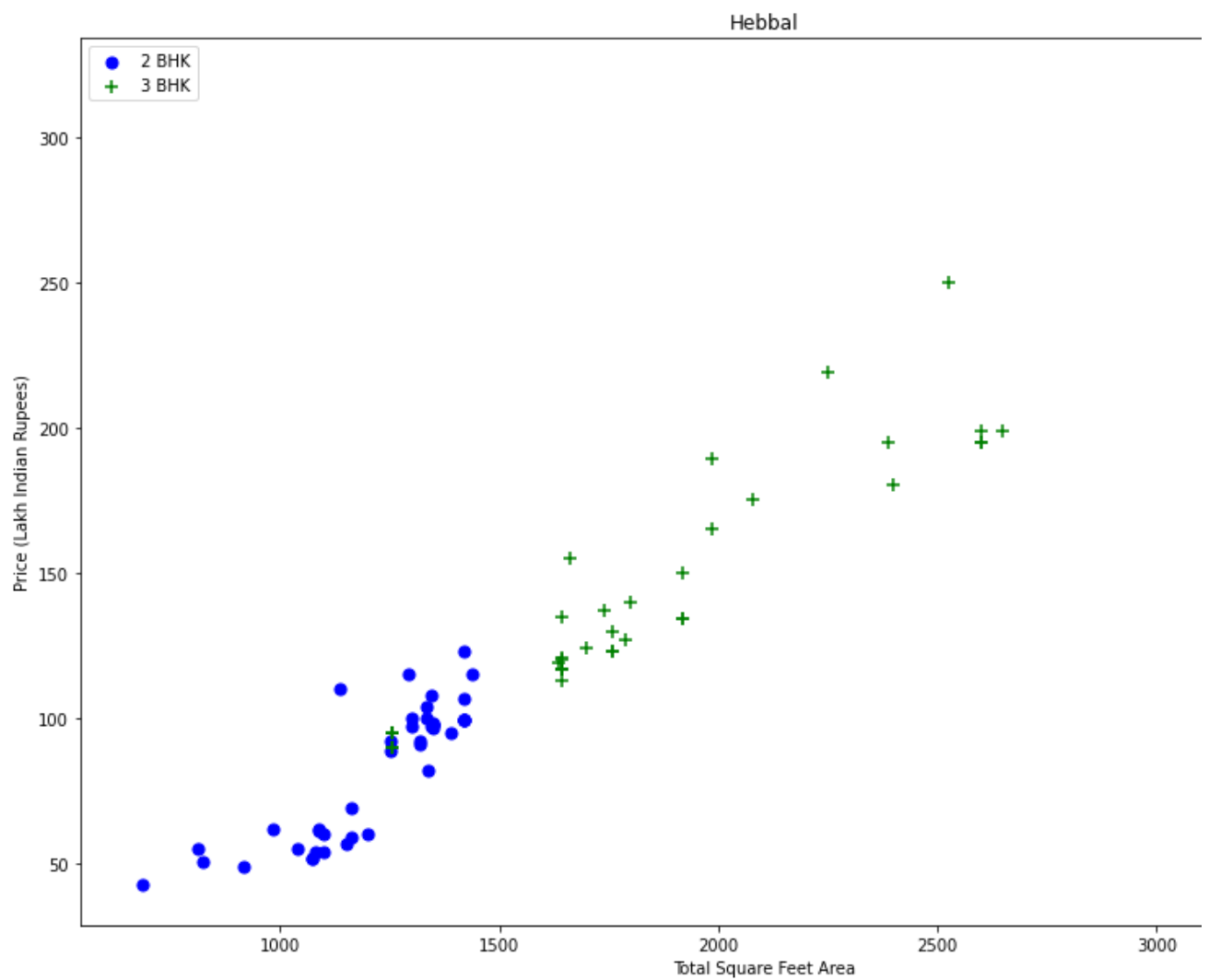
↳ (7317, 7)

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

↳



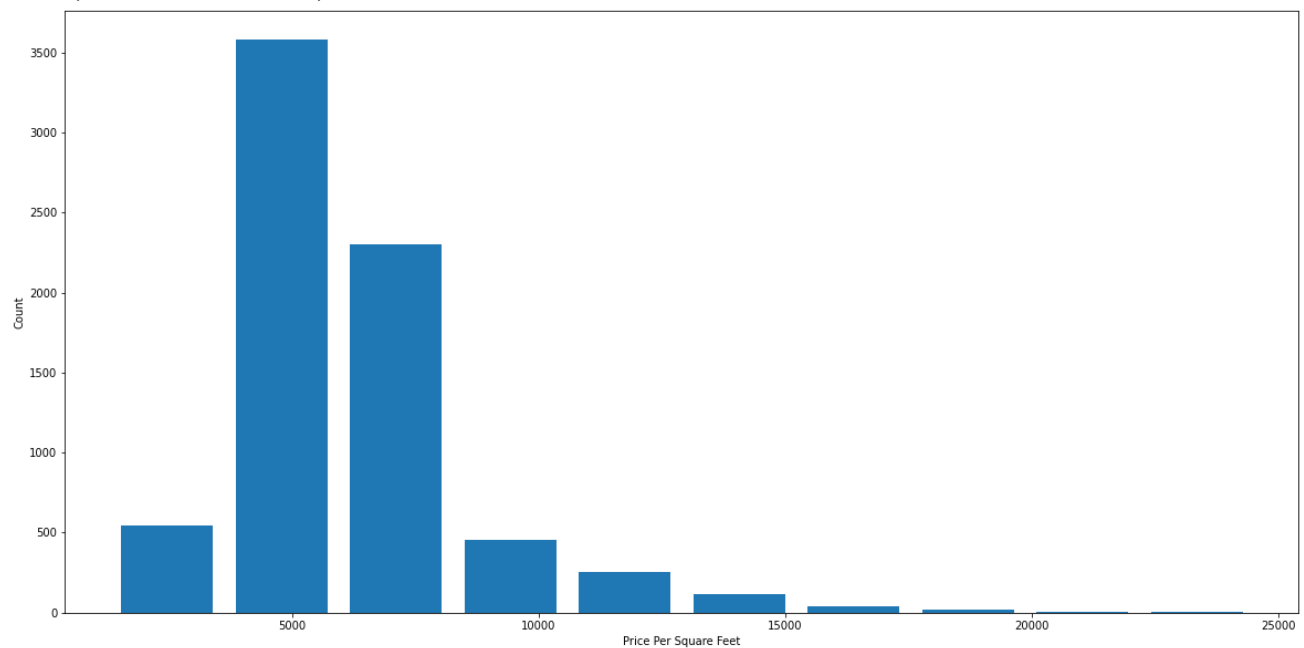
```
plot_scatter_chart(df8,"Hebbal")
```



```
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")
```

```
↳ Text(0, 0.5, 'Count')
```



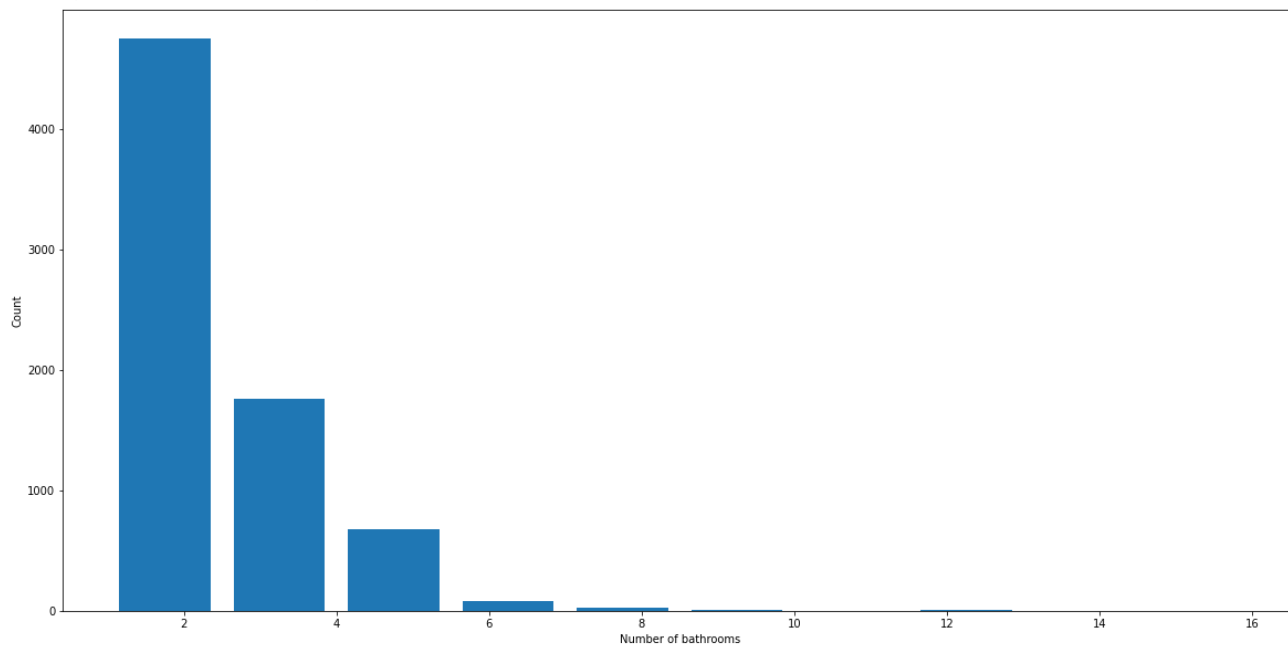
```
df8.bath.unique()
```

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

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

```
↳
```

Text(0, 0.5, 'Count')



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 |
| 8483 | other | 10 BHK | 12000.0 | 12.0 | 525.0 | 10 | 4375.000000 |
| 8572 | other | 16 BHK | 10000.0 | 16.0 | 550.0 | 16 | 5500.000000 |
| 9306 | other | 11 BHK | 6000.0 | 12.0 | 150.0 | 11 | 2500.000000 |
| 9637 | other | 13 BHK | 5425.0 | 13.0 | 275.0 | 13 | 5069.124424 |

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 |
| 8408 | other | 6 BHK | 11338.0 | 9.0 | 1000.0 | 6 | 8819.897689 |

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

df9.shape



(7239, 7)

df9.head(2)



| | 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 |

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



| | 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 |

```
dummies = pd.get_dummies(df10.location)
dummies.head(3)
```



| | 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 rows × 241 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 | L |
|---|---------------------|------------|------|-------|-----|---------------------|--------------------|---------------------------|----------------------|---|
| 0 | 1st Block Jayanagar | 2850.0 | 4.0 | 428.0 | 4 | 1 | 0 | 0 | 0 | |
| 1 | 1st Block Jayanagar | 1630.0 | 3.0 | 194.0 | 3 | 1 | 0 | 0 | 0 | |
| 2 | 1st Block Jayanagar | 1875.0 | 2.0 | 235.0 | 3 | 1 | 0 | 0 | 0 | |
| 3 | 1st Block Jayanagar | 1200.0 | 2.0 | 130.0 | 3 | 1 | 0 | 0 | 0 | |
| 4 | 1st Block Jayanagar | 1235.0 | 2.0 | 148.0 | 2 | 1 | 0 | 0 | 0 | |

5 rows × 245 columns

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

↗

| | 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 | |
| 1 | 1630.0 | 3.0 | 194.0 | 3 | 1 | 0 | 0 | 0 | 0 | |

2 rows × 244 columns

```
df12.shape
```

↗

| |
|-------------|
| (7239, 244) |
|-------------|

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

↗

| | 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 | |
| 1 | 1630.0 | 3.0 | 3 | 1 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 1875.0 | 2.0 | 3 | 1 | 0 | 0 | 0 | 0 | 0 | |

3 rows × 243 columns

```
df1 = pd.read_csv("/content/drive/My Drive/datasets_20710_26737_Bengaluru_House_Data.csv")
df1.head()
```

↗

| | area_type | availability | location | size | society | total_sqft | bath | ba |
|---|---------------------|---------------|--------------------------|-----------|---------|------------|------|----|
| 0 | Super built-up Area | 19-Dec | Electronic City Phase II | 2 BHK | Coomee | 1056 | 2.0 | |
| 1 | Plot Area | Ready To Move | Chikka Tirupathi | 4 Bedroom | Theanmp | 2600 | 5.0 | |
| 2 | Built-up Area | Ready To Move | Uttarahalli | 3 BHK | NaN | 1440 | 2.0 | |

```
df1.replace("Super built-up Area","cool")
```



| | area_type | availability | location | size | society | total_sqft | bath |
|-------|---------------------|---------------|--------------------------|-----------|---------|------------|------|
| 0 | Super built-up Area | 19-Dec | Electronic City Phase II | 2 BHK | Coomee | 1056 | 2.0 |
| 1 | Plot Area | Ready To Move | Chikka Tirupathi | 4 Bedroom | Theanmp | 2600 | 5.0 |
| 2 | Built-up Area | Ready To Move | Uttarahalli | 3 BHK | NaN | 1440 | 2.0 |
| 3 | Super built-up Area | Ready To Move | Lingadheeranahalli | 3 BHK | Soiewre | 1521 | 3.0 |
| 4 | Super built-up Area | Ready To Move | Kothanur | 2 BHK | NaN | 1200 | 2.0 |
| ... | ... | ... | ... | ... | ... | ... | .. |
| 13315 | Built-up Area | Ready To Move | Whitefield | 5 Bedroom | ArsiaEx | 3453 | 4.0 |