

```
In [192]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import matplotlib
matplotlib.rcParams["figure.figsize"]=(20,10)
```

```
In [193]: df1=pd.read_csv("C:/Users/HP/Desktop/Bengaluru_House_Data.csv")
df1.head()
```

Out[193]:

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

Out[194]: (13320, 5)

```
In [195]: df2.isnull().sum()
```

```
Out[195]: location      1
          size         16
          total_sqft    0
          bath         73
          price         0
          dtype: int64
```

```
In [196]: df2.shape
```

```
Out[196]: (13320, 5)
```

```
In [197]: df3 = df2.dropna()
          df3.isnull().sum()
```

```
Out[197]: location      0
          size          0
          total_sqft    0
          bath          0
          price         0
          dtype: int64
```

```
In [198]: df3.shape
```

```
Out[198]: (13246, 5)
```

```
In [199]: df3['bhk'] = df3['size'].apply(lambda x: int(x.split(' ')[0]))
          df3.bhk.unique()
```

```
C:\Users\HP\anaconda3\lib\site-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/user_guide/indexing.html#returning-a-view-versus-a-copy
"""Entry point for launching an IPython kernel.
```

```
Out[199]: array([ 2,  4,  3,  6,  1,  8,  7,  5, 11,  9, 27, 10, 19, 16, 43, 14,
```

```
12,  
    13, 18], dtype=int64)
```

```
In [200]: def is_float(x):  
          try:  
            float(x)  
          except:  
            return False  
          return True
```

```
In [201]: df3[~df3['total_sqft'].apply(is_float)].head(10)
```

Out[201]:

	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

```
In [202]: 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 [203]: df4 = df3.copy()
df4.total_sqft = df4.total_sqft.apply(convert_sqft_to_num)
df4 = df4[df4.total_sqft.notnull()]
df4.head(2)
```

```
Out[203]:
```

	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

```
In [204]: df4.loc[30]
```

```
Out[204]: location      Yelahanka
size              4 BHK
total_sqft        2475
bath              4
price             186
bkh              4
Name: 30, dtype: object
```

```
In [205]: (2100+2850)/2
```

```
Out[205]: 2475.0
```

```
In [206]: df5 = df4.copy()
df5['price_per_sqft'] = df5['price']*100000/df5['total_sqft']
df5.head()
```

```
Out[206]:
```

	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

	location	size	total_sqft	bath	price	bhk	price_per_sqft
4	Kothanur	2 BHK	1200.0	2.0	51.00	2	4250.000000

```
In [207]: df5_stats = df5['price_per_sqft'].describe()
df5_stats
```

```
Out[207]: 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
```

```
In [208]: df5.to_csv("C:/Users/HP/Desktop/bhp.csv", index=False)
```

```
In [209]: df5.location = df5.location.apply(lambda x: x.strip())
location_stats = df5['location'].value_counts(ascending=False)
location_stats
```

```
Out[209]: Whitefield          533
Sarjapur Road          392
Electronic City        304
Kanakpura Road         264
Thanisandra            235
...
Rajarajesheari nagar    1
Attiguppe               1
Hoskote near            1
5th block Koramangala   1
Chambenahalli           1
Name: location, Length: 1287, dtype: int64
```

```
In [210]: location_stats.values.sum()
```

Out[210]: 13200

```
In [211]: len(location_stats[location_stats>10])
```

Out[211]: 240

```
In [212]: len(location_stats)
```

Out[212]: 1287

```
In [213]: len(location_stats[location_stats<=10])
```

Out[213]: 1047

```
In [214]: location_stats_less_than_10 = location_stats[location_stats<=10]
location_stats_less_than_10
```

Out[214]:

Thyagaraja Nagar	10
Nagadevanahalli	10
Basapura	10
Sadashiva Nagar	10
Nagappa Reddy Layout	10
..	
Rajarajesheari nagar	1
Attiguppe	1
Hoskote near	1
5th block Koramangala	1
Chambenahalli	1
Name: location, Length: 1047, dtype: int64	

```
In [215]: len(df5.location.unique())
```

Out[215]: 1287

```
In [216]: df5.location = df5.location.apply(lambda x: 'other' if x in location_stats_less_than_10 else x)
len(df5.location.unique())
```

Out[216]: 241

In [217]: `df5.head(10)`

Out[217]:

	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

In [218]: `df5[df5.total_sqft/df5.bhk<300].head()`

Out[218]:

	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

In [219]: `df5.shape`

Out[219]: (13200, 7)

```
In [220]: df6 = df5[~(df5.total_sqft/df5.bhk<300)]
df6.shape
```

```
Out[220]: (12456, 7)
```

```
In [221]: df6.price_per_sqft.describe()
```

```
Out[221]: 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
```

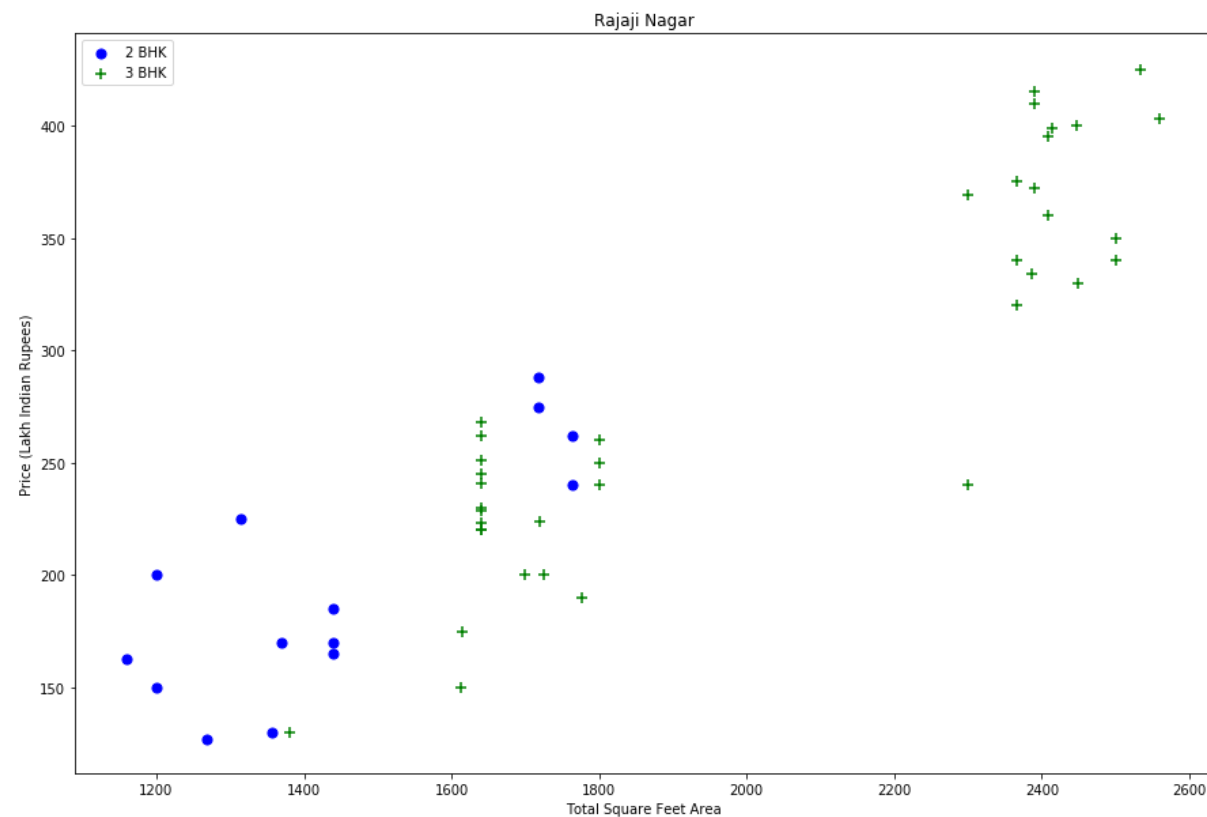
```
In [222]: def remove_pps_outliers(df):
df_out = pd.DataFrame()
for key, subdf in df.groupby('location'):
    m = np.mean(subdf.price_per_sqft)
    st = np.std(subdf.price_per_sqft)
    reduced_df = subdf[(subdf.price_per_sqft>(m-st)) & (subdf.price_per_sqft<=(m+st))]
    df_out = pd.concat([df_out,reduced_df],ignore_index=True)
return df_out
df7 = remove_pps_outliers(df6)
df7.shape
```

```
Out[222]: (10242, 7)
```

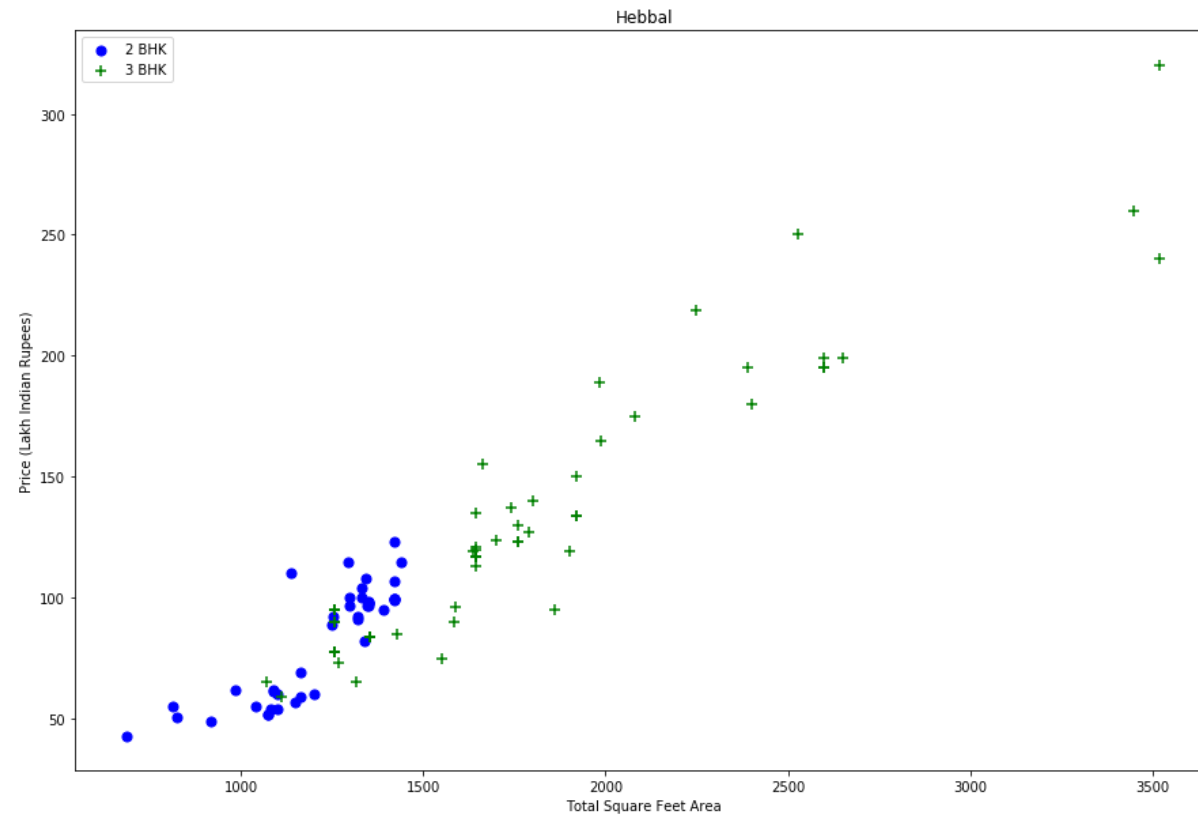
```
In [223]: def plot_scatter_chart(df,location):
bhk2 = df[(df.location==location) & (df.bhk==2)]
bhk3 = df[(df.location==location) & (df.bhk==3)]
matplotlib.rcParams['figure.figsize'] = (15,10)
plt.scatter(bhk2.total_sqft,bhk2.price,color='blue',label='2 BHK',
s=50)
plt.scatter(bhk3.total_sqft,bhk3.price,marker='+', color='green',la
```



```
bel='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")
```



```
In [224]: plot_scatter_chart(df7,"Hebbal")
```

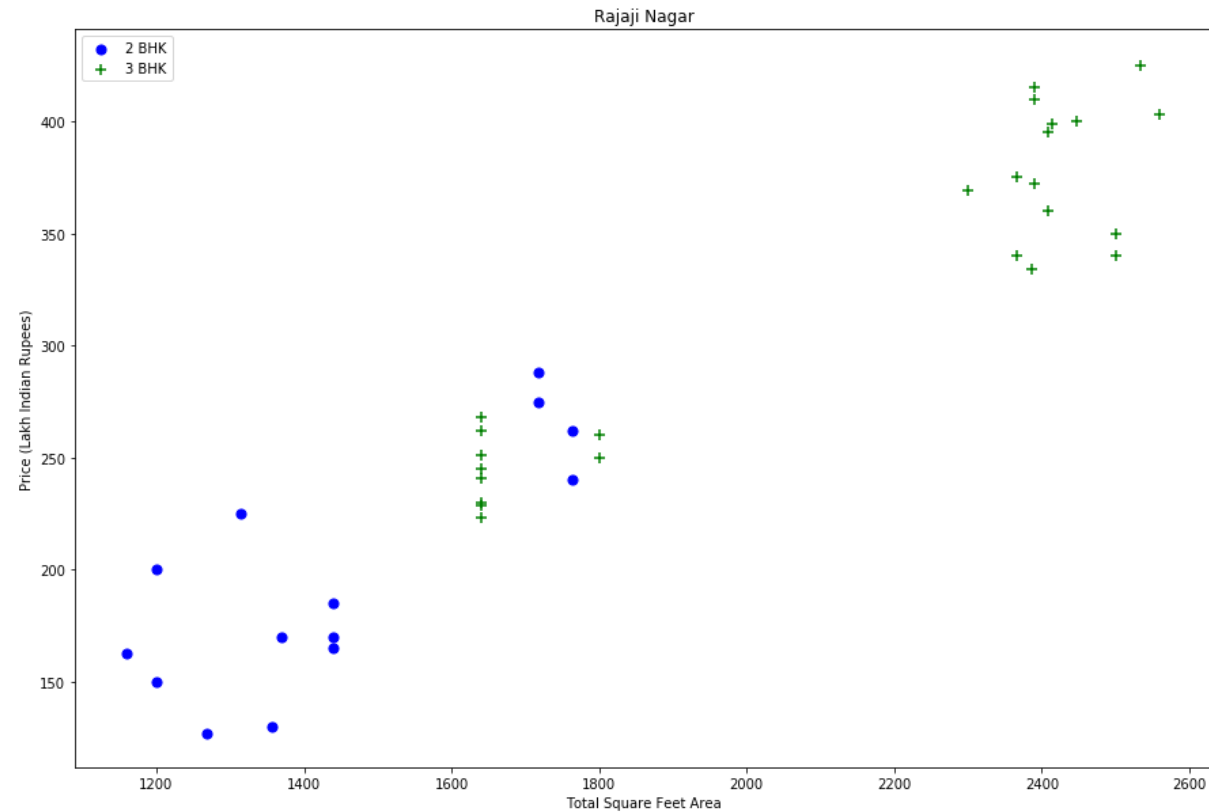


```
In [225]: def remove_bhk_outliers(df):
            exclude_indices = np.array([])
            for location, location_df in df.groupby('location'):
                bhk_stats = {}
                for bhk, bhk_df in location_df.groupby('bhk'):
                    bhk_stats[bhk] = {
                        '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
```

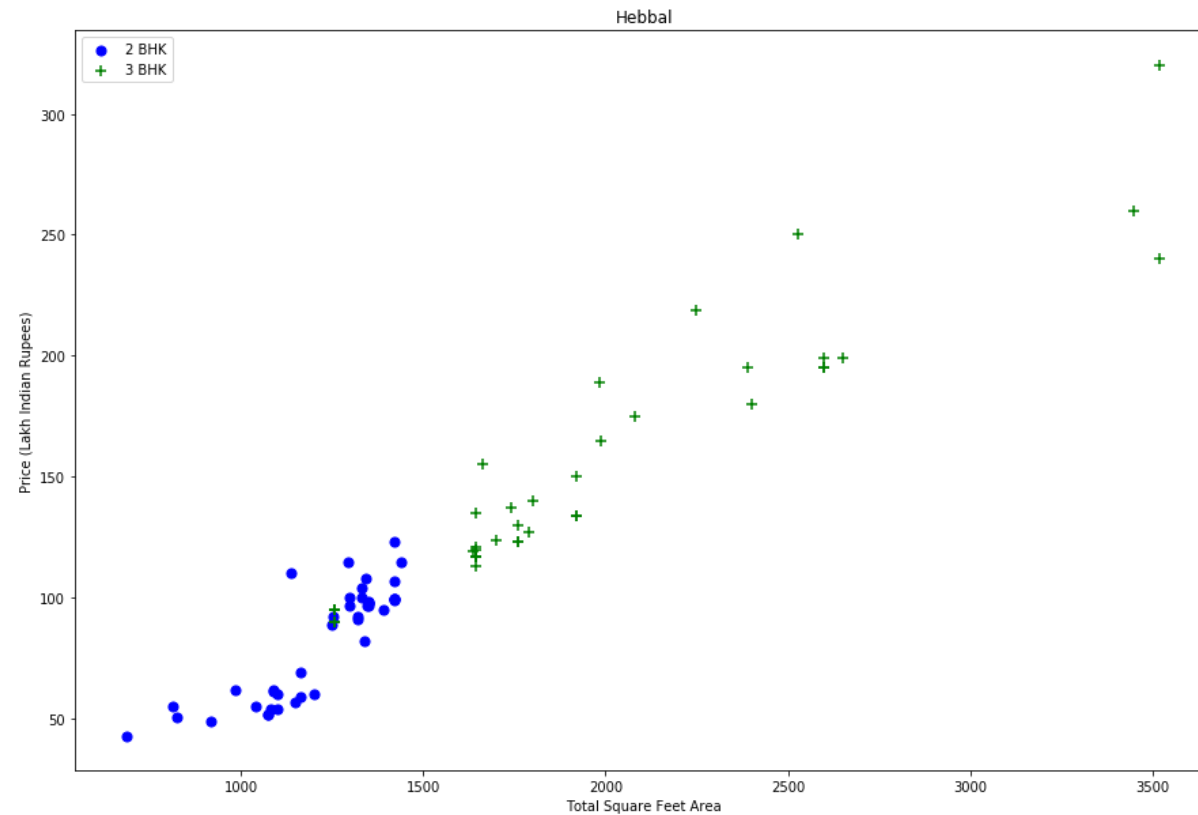
```
_df.price_per_sqft<(stats['mean'])).index.values)
    return df.drop(exclude_indices,axis='index')
df8 = remove_bhk_outliers(df7)
# df8 = df7.copy()
df8.shape
```

Out[225]: (7317, 7)

In [226]: plot_scatter_chart(df8,"Rajaji Nagar")

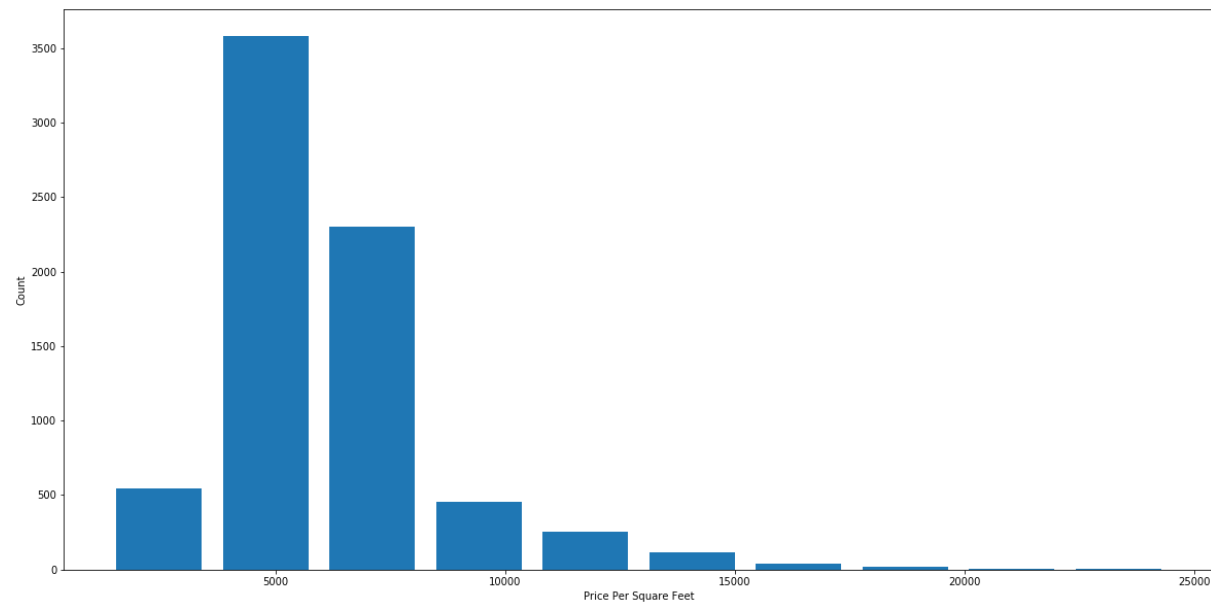


In [227]: plot_scatter_chart(df8,"Hebbal")



```
In [228]: 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[228]: Text(0, 0.5, 'Count')
```

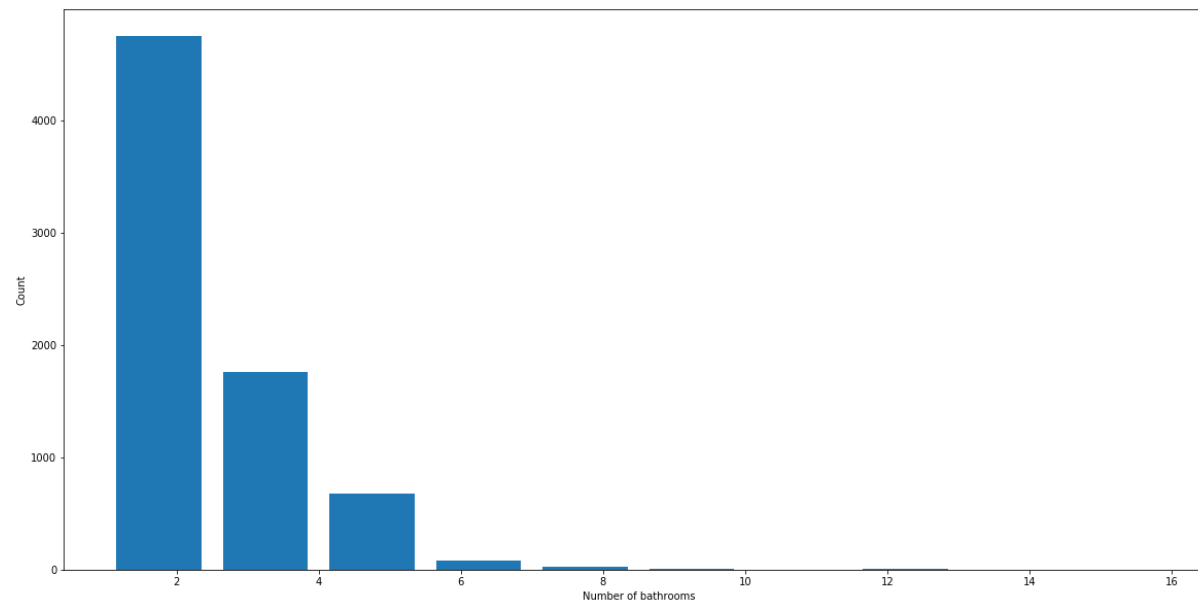


```
In [229]: df8.bath.unique()
```

```
Out[229]: array([ 4.,  3.,  2.,  5.,  8.,  1.,  6.,  7.,  9., 12., 16., 13.] )
```

```
In [230]: plt.hist(df8.bath,rwidth=0.8)
plt.xlabel("Number of bathrooms")
plt.ylabel("Count")
```

```
Out[230]: Text(0, 0.5, 'Count')
```



```
In [231]: df8[df8.bath>10]
```

Out[231]:

	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

```
In [232]: df8[df8.bath>df8.bhk+2]
```

Out[232]:

	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

	location	size	total_sqft	bath	price	bhk	price_per_sqft
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

```
In [233]: df9 = df8[df8.bath<df8.bhk+2]
df9.shape
```

Out[233]: (7239, 7)

```
In [234]: df9.head(2)
```

Out[234]:

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

Out[235]:

	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

```
In [236]: dummies = pd.get_dummies(df10.location)
dummies.head(3)
```

Out[236]:

	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	...	Vis
0	1	0	0	0	0	0	0	0	0	0	...	

	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	...	Vis
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

◀		▶
---	--	---

In [237]: `df11 = pd.concat([df10,dummies.drop('other',axis='columns')],axis='columns')`
`df11.head()`

Out[237]:

	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 × 245 columns

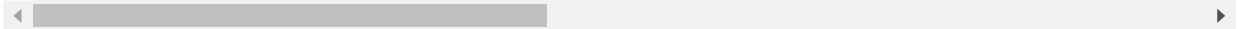
◀		▶
---	--	---

In [238]: `df12 = df11.drop('location',axis='columns')`
`df12.head(2)`

Out[238]:

	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	...	Vija
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 × 244 columns



In [239]: df12.shape

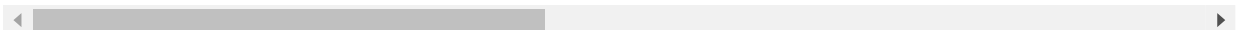
Out[239]: (7239, 244)

In [240]: X = df12.drop(['price'],axis='columns')
X.head(3)

Out[240]:

	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	...	Vija
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 × 243 columns



In [241]: X.shape

Out[241]: (7239, 243)

In [242]: y = df12.price
y.head(3)

```
Out[242]: 0    428.0  
          1    194.0  
          2    235.0  
          Name: price, dtype: float64
```

```
In [243]: len(y)
```

```
Out[243]: 7239
```

```
In [244]: 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,r  
andom_state=10)
```

```
In [245]: from sklearn.linear_model import LinearRegression  
lr_clf = LinearRegression()  
lr_clf.fit(X_train,y_train)  
lr_clf.score(X_test,y_test)
```

```
Out[245]: 0.8629132245229442
```

```
In [246]: 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[246]: array([0.82702546, 0.86027005, 0.85322178, 0.8436466 , 0.85481502])
```

```
In [247]: 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': {
            'normalize': [True, False]
        },
    },
    '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=0.2, random_state=0)
for algo_name, config in algos.items():
    gs = GridSearchCV(config['model'], config['params'], cv=cv, re
turn_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)

```

Out[247]:

	model	best_score	best_params
0	linear_regression	0.847796	{'normalize': False}

	model	best_score	best_params
1	lasso	0.726774	{'alpha': 2, 'selection': 'random'}
2	decision_tree	0.717352	{'criterion': 'mse', 'splitter': 'best'}

```
In [131]: 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]
```

location , sqr feet , bedroom , bathroom

```
In [132]: predict_price('1st Phase JP Nagar',1000, 2, 2)
```

```
Out[132]: 83.86570258312275
```

```
In [198]: predict_price('1st Phase JP Nagar',1000, 3, 3)
```

```
Out[198]: 86.08062284987054
```

```
In [199]: predict_price('Indira Nagar',1000, 2, 2)
```

```
Out[199]: 193.31197733179937
```

This is how we can predict the prices MADE BY ATUL KOHAR

In []:

Export the tested model to a pickle file

```
In [200]: import pickle
with open('bangalore_home_prices_model.pickle','wb') as f:
    pickle.dump(lr_clf,f)
```

In []:

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

```
-----
NameError                                Traceback (most recent call l
ast)
<ipython-input-1-32a87433566e> in <module>
      1 import json
      2 column = {
----> 3     'data_columns': [col.lower() for col in X.columns]
      4 }
      5 with open("columns.json", "w") as f:

NameError: name 'X' is not defined
```

In []: