

# price-prediction-1

December 4, 2023

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

```
[3]: df1 = pd.read_csv("Bengaluru_House_Data.csv")
df1.head()
```

```
[3]:
```

		area_type	availability	location	size \
0	Super built-up	Area	19-Dec	Electronic City Phase II	2 BHK
1	Plot	Area	Ready To Move	Chikka Tirupathi	4 Bedroom
2	Built-up	Area	Ready To Move	Uttarahalli	3 BHK
3	Super built-up	Area	Ready To Move	Lingadheeranahalli	3 BHK
4	Super built-up	Area	Ready To Move	Kothanur	2 BHK

	society	total_sqft	bath	balcony	price
0	Coomee	1056	2.0	1.0	39.07
1	Theanmp	2600	5.0	3.0	120.00
2	NaN	1440	2.0	3.0	62.00
3	Soiewre	1521	3.0	1.0	95.00
4	NaN	1200	2.0	1.0	51.00

```
[4]: df1.shape
```

```
[4]: (13320, 9)
```

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

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

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

```
[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

```
[7]: df2.isnull().sum()
```

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

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

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

```
[9]: df3.shape
```

```
[9]: (13246, 5)
```

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

```
[10]: 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)
```

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

```
C:\Users\admin\AppData\Local\Temp\ipykernel_13200\2222900254.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](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]))
```

```
[12]: df3.head()
```

```
[12]:
```

	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

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

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

```
[14]: df3[df3.bhk>20]
```

```
[14]:
```

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

```
[15]: df3.total_sqft.unique()
```

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

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

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

```
[17]:
```

	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

```
[18]: 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
```

```
[19]: convert_sqft_to_num('2100 - 2850')
```

```
[19]: 2475.0
```

```
[20]: df4 = df3.copy()
      df4['total_sqft'] = df4['total_sqft'].apply(convert_sqft_to_num)
      df4.head(3)
```

	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

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

	price_per_sqft
0	3699.810606
1	4615.384615
2	4305.555556
3	6245.890861
4	4250.000000

```
[22]: len(df5.location.unique())
```

```
[22]: 1304
```

```
[23]: df5.location = df5.location.apply(lambda x: x.strip())
location_stats = df5.groupby('location')['location'].agg('count').
    ↪sort_values(ascending = False)
location_stats
```

```
[23]: 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
```

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

```
[24]: 1052
```

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

```
[25]: 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
```

```
[26]: len(df5.location.unique())
```

```
[26]: 1293
```

```
[27]: df5.location = df5.location.apply(lambda x: 'other' if x in_
    ↪location_stats_less_than_10 else x)
len(df5.location.unique())
```

[27]: 242

```
[28]: df5.head()
```

```
[28]:
```

	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	

	price_per_sqft
0	3699.810606
1	4615.384615
2	4305.555556
3	6245.890861
4	4250.000000

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

```
[29]:
```

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

	price_per_sqft
9	36274.509804
45	33333.333333
58	10660.980810
68	6296.296296
70	20000.000000

```
[30]: df5.shape
```

[30]: (13246, 7)

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

[31]: (12502, 7)

```
[32]: df6.price_per_sqft.describe()
```

```
[32]: 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

```

```

[33]: 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

```

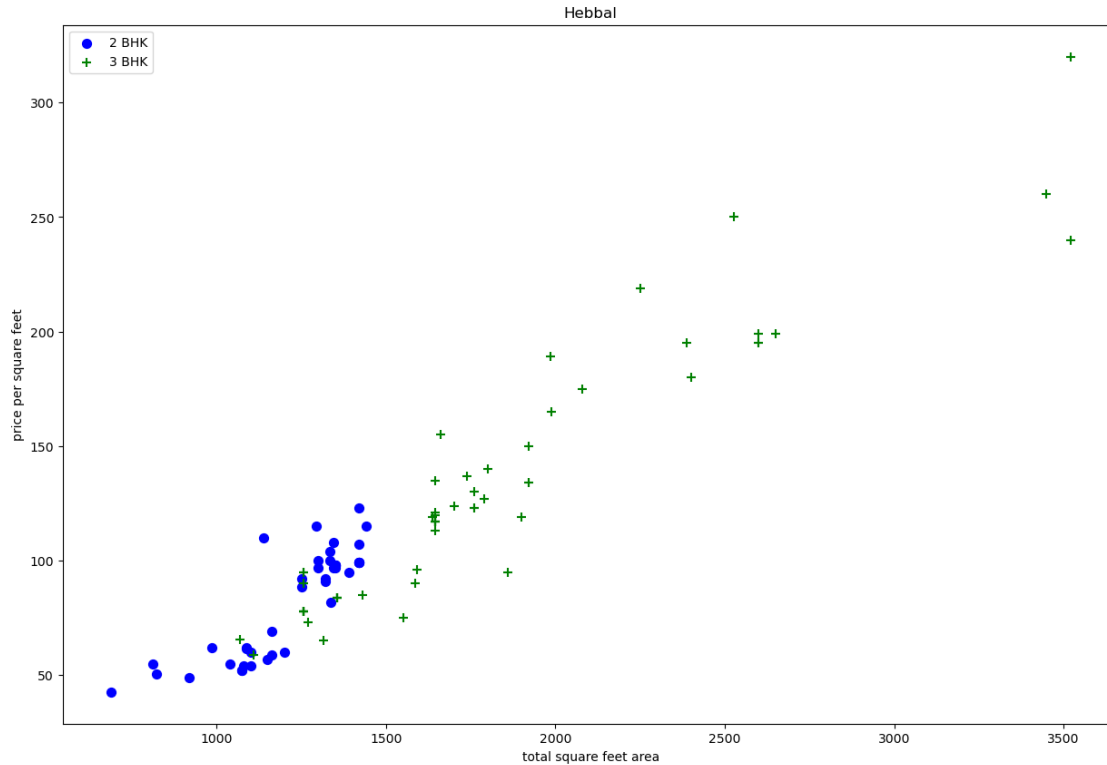
```
[33]: (10241, 7)
```

```

[34]: 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 per square feet")
        plt.title(location)
        plt.legend()

plot_scatter_chart(df7,"Hebbal")

```



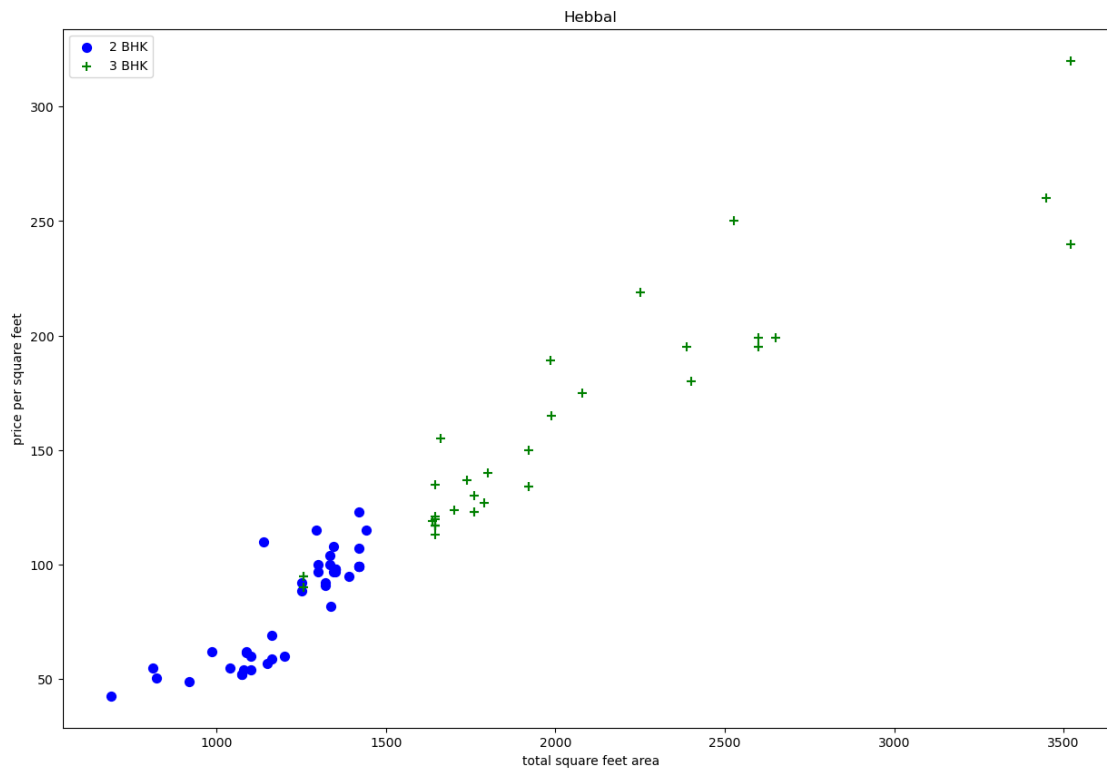
```
[35]: def remove_bhk_outliers(df):
    exclude_indices = np.array([])
    for location, location_df in df.groupby('location'):
        bhk_stats= {}
        for bhk, bhk_df in location_df.groupby('bhk'):
            bhk_stats[bhk] = {
                'mean': np.mean(bhk_df.price_per_sqft),
     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

[35]: (7329, 7)

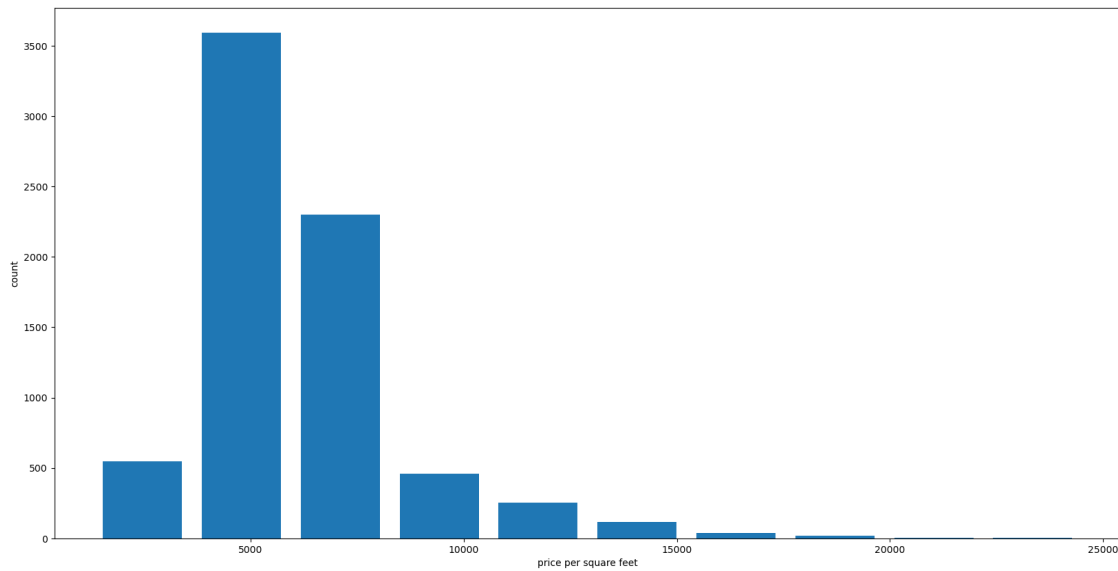


```
[36]: plot_scatter_chart(df8, "Hebbal")
```



```
[37]: 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")
```

```
[37]: Text(0, 0.5, 'count')
```



```
[38]: df8.bath.unique()
```

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

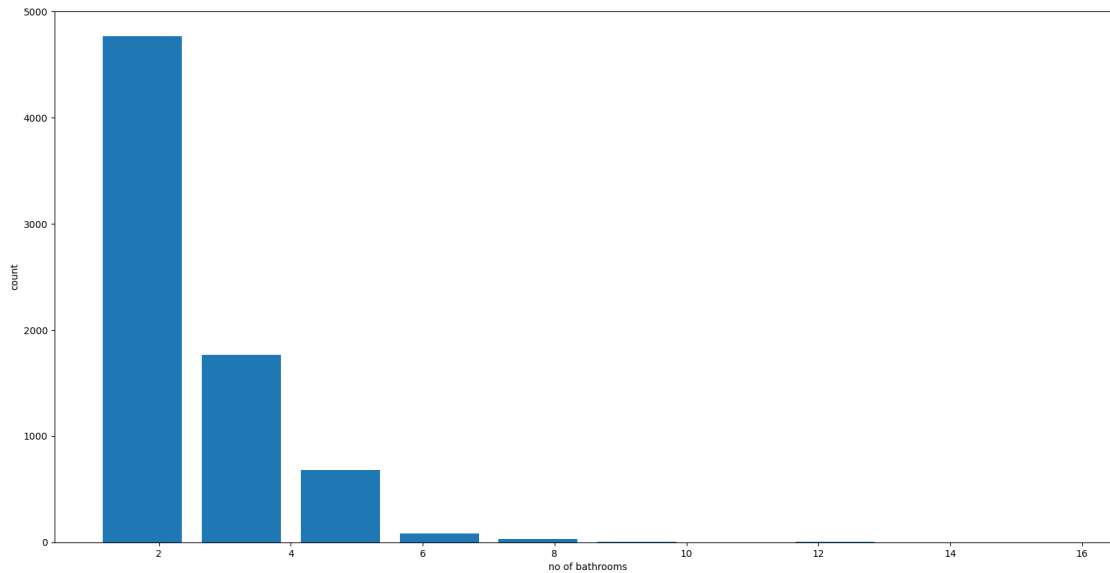
```
[39]: df8[df8.bath>10]
```

```
[39]:
```

	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

```
[40]: plt.hist(df8.bath,rwidth=0.8)
plt.xlabel("no of bathrooms")
plt.ylabel("count")
```

```
[40]: Text(0, 0.5, 'count')
```



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

```
[41]:
```

	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

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

```
[42]: (7251, 7)
```

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

```
[43]:
```

	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

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

```
[45]:
```

	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	\
0	1	0	0	
1	1	0	0	
2	1	0	0	

	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	\
0	0	0	0	
1	0	0	0	
2	0	0	0	

	6th Phase JP Nagar	7th Phase JP Nagar	8th Phase JP Nagar	\
0	0	0	0	
1	0	0	0	
2	0	0	0	

	9th Phase JP Nagar	...	Vishveshwarya Layout	Vishwapriya Layout	\
0	0	...	0	0	
1	0	...	0	0	
2	0	...	0	0	

	Vittasandra	Whitefield	Yelachenahalli	Yelahanka	Yelahanka New Town	\
0	0	0	0	0	0	
1	0	0	0	0	0	
2	0	0	0	0	0	

	Yelenahalli	Yeshwanthpur	other
0	0	0	0
1	0	0	0
2	0	0	0

[3 rows x 242 columns]

```
[46]: df11 = pd.concat([df10,dummies.drop('other',axis='columns')],axis='columns')
df11.head(3)
```

```
[46]:
```

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

	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	\
0	0	0	0	
1	0	0	0	
2	0	0	0	

	5th Block Hbr Layout	...	Vijayanagar	Vishveshwarya Layout	\
0	0	...	0	0	
1	0	...	0	0	
2	0	...	0	0	

	Vishwapriya Layout	Vittasandra	Whitefield	Yelachenahalli	Yelahanka	\
--	--------------------	-------------	------------	----------------	-----------	---

0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0

	Yelahanka New Town	Yelenahalli	Yeshwanthpur
0	0	0	0
1	0	0	0
2	0	0	0

[3 rows x 246 columns]

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

```
[47]: total_sqft  bath  price  bhk  1st Block Jayanagar  1st Phase JP Nagar \
0      2850.0    4.0  428.0   4                1                0
1      1630.0    3.0  194.0   3                1                0

2nd Phase Judicial Layout  2nd Stage Nagarbhavi  5th Block Hbr Layout \
0                        0                        0                0
1                        0                        0                0

5th Phase JP Nagar  ...  Vijayanagar  Vishveshwarya Layout \
0                0  ...                0                0
1                0  ...                0                0

Vishwapriya Layout  Vittasandra  Whitefield  Yelachenahalli  Yelahanka \
0                0                0                0                0
1                0                0                0                0

Yelahanka New Town  Yelenahalli  Yeshwanthpur
0                0                0                0
1                0                0                0
```

[2 rows x 245 columns]

```
[48]: df12.shape
```

```
[48]: (7251, 245)
```

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

```
[50]: total_sqft  bath  bhk  1st Block Jayanagar  1st Phase JP Nagar \
0      2850.0    4.0   4                1                0
1      1630.0    3.0   3                1                0
2      1875.0    2.0   3                1                0
```

3	1200.0	2.0	3	1	0
4	1235.0	2.0	2	1	0

	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	

	5th Phase JP Nagar	6th Phase JP Nagar	... Vijayanagar	\
0	0	0	...	0
1	0	0	...	0
2	0	0	...	0
3	0	0	...	0
4	0	0	...	0

	Vishveshwarya Layout	Vishwapriya Layout	Vittasandra	Whitefield	\
0	0	0	0	0	
1	0	0	0	0	
2	0	0	0	0	
3	0	0	0	0	
4	0	0	0	0	

	Yelachenahalli	Yelahanka	Yelahanka New Town	Yelenahalli	Yeshwanthpur
0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0

[5 rows x 244 columns]

```
[51]: y = df12.price
      y.head()
```

```
[51]: 0    428.0
      1    194.0
      2    235.0
      3    130.0
      4    148.0
      Name: price, dtype: float64
```

```
[53]: 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=10)
```

```
[55]: from sklearn.linear_model import LinearRegression
lr_clf = LinearRegression()
lr_clf.fit(X_train,y_train)
```

```
[55]: LinearRegression()
```

```
[56]: lr_clf.score(X_test,y_test)
```

```
[56]: 0.8452277697874279
```

```
[60]: 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)
```

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

```
[63]: 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]
```

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

D:\anaconda\Lib\site-packages\sklearn\base.py:464: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names  
warnings.warn(

```
[64]: 83.49904677167738
```

```
[65]: predict_price('Vijayanagar',1000,2,2)
```

D:\anaconda\Lib\site-packages\sklearn\base.py:464: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names  
warnings.warn(

```
[65]: 62.29709374581037
```

```
[66]: import pickle
      with open('house_price_prediction_bangalore','wb') as f:
          pickle.dump(lr_clf,f)
```

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