

IMPORT ALL REQUIRED LIBRARIES

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

LOAD DATASET

```
df = pd.read_csv('Bengaluru_House_Data.csv')
df.head()
```

	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 Trupathi	4 Bedroom	Theamp	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	Solewa	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

```
df.shape
```

```
(13328, 9)
```

```
df.isnull().mean()*100
```

area_type	0.000000
availability	0.000000
location	0.007508
size	0.128128
society	41.398386
total_sqft	0.000000
bath	0.546898
balcony	4.572072
price	0.000000
dtype:	float64

```
df['area_type'].value_counts()
```

Super built-up Area	8796
Built-up Area	2418
Plot Area	2025
Carpet Area	87
Name: area_type, dtype: int64	

```
df.drop(columns=['availability','area_type','society','balcony'],axis=1,inplace=True)
```

```
df.head()
```

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Trupathi	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

```
df.isnull().sum()
```

location	1
size	18
total_sqft	0
bath	73
price	0
dtype:	int64

```
df.dropna(inplace=True)
```

```
df.isnull().sum()
```

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

```
df.shape
```

```
(13246, 5)
```

```
df.head()
```

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Trupathi	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

```
df['size'].unique()
```

```
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', '10 BHK', '16 BHK', '43 Bedroom', '14 BHK', '8 BHK', '12 Bedroom', '13 BHK', '18 Bedroom'], dtype=object)
```

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

```
df.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
1	Chikka Trupathi	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

```
df.drop(columns=["size"],axis=1,inplace=True)
```

```
df.shape
```

```
(13246, 5)
```

```
df[df.bhk==2]
```

	location	total_sqft	bath	price	bhk
1718	Electronic City Phase II	8000	27.0	230.0	27
4684	Munnekotai	2400	40.0	660.0	43

```
df.total_sqft.isnull()
```

```
array(['1056', '12600', '1440', ..., '1133', '1384', '774', '4680'], dtype=object)
```

```
def is_float(x):
```

```
    try:
        float(x)
    except:
        return False
    return True
```

```
df[df['total_sqft'].apply(is_float)].head(10)
```

	location	total_sqft	bath	price	bhk
30	Yelahanka	2100-2850	4.0	186.00	4
122	Hebbal	3067-8156	4.0	477.000	4
137	8th Phase JP Nagar	1042-1105	2.0	54.005	2
105	Sarjapur	1145-1340	2.0	43.490	2
188	KR Puram	1015-1540	2.0	56.800	2
410	Kengeri	34.46Sq Meter	1.0	18.500	1
549	Hennur Road	1195-1440	2.0	63.770	2
648	Anekeere	4125Perch	9.0	265.000	9
661	Yelahanka	1120-1145	2.0	48.130	2
672	Bottahaloor	3090-5002	4.0	445.000	4

```
def convert_sqft_into_number(x):
    token = x.split('-')
    if len(token) == 2:
        return (float(token[0]) + float(token[1])) / 2
    try:
        return float(x)
    except:
        return None
```

```
df1 = df.copy()
```

```
df1['total_sqft'] = df1['total_sqft'].apply(convert_sqft_into_number)
```

```
df1.loc[30]
```

location	Yelahanka
total_sqft	2475.0
bath	4.0
price	186.0
bhk	4
Name: 30, dtype: object	

```
df2 = df1.copy()
```

```
df2['price_per_sqft'] = df2['price']/100000 / df2['total_sqft']
```

```
df2.head()
```

	location	total_sqft	bath	price	bhk	price_per_sqft
0	Electronic City Phase II	1056.0	2.0	39.07	2	3699.810066
1	Chikka Trupathi	2600.0	5.0	120.00	4	4615.384615
2	Uttarahalli	1440.0	2.0	62.00	3	4305.555556
3	Lingadheeranahalli	1521.0	3.0	95.00	3	6245.890861
4	Kothanur	1200.0	2.0	51.00	2	4250.000000

```
df2['location'].value_counts()
```

Whitefield	534
Sarjapur Road	392
Electronic City	302
Kanakapura Road	266
Thansandra	233

...	1
Vidyapeeta	1
Maruthi Extension	1
Okaligura	1
Old Town	1
Abshot Layout	1
Name: location, Length: 1384, dtype: int64	

```
df2['location'] = df2['location'].apply(lambda x: x.strip())
```

```
df2.location.value_counts()
```

Whitefield	535
Sarjapur Road	392
Electronic City	304
Kanakapura Road	268
Thansandra	236

Vasantapura main road	1
Bapuji Layout	1
1st Stage Nadra Krishna Layout	1
BEML Layout 5th stage	1
Abshot Layout	1
Name: location, Length: 1293, dtype: int64	

```
location[location==10]
```

```
loc_less_than_10 = loc_stats[loc_stats<=10]
```

```
loc_less_than_10
```

```
df2.location = df2.location.apply(lambda x: 'other' if x in loc_less_than_10 else x)
```

```
df2.head()
```

```
len(df2.location.unique())
```

```
1293
```

```
df2[(df2.total_sqft / df2.bhk < 300)].head()
```

	location	total_sqft	bath	price	bhk	price_per_sqft
9	Gandhi Bazar	1020.0	6.0	370.0	6	36274.509804
45	HSR Layout	600.0	9.0	200.0	8	33333.333333
58	Munageshpalya	1407.0	4.0	150.0	6	10660.980810
68	Devarachikkanahalli	1350.0	7.0	85.0	8	6296.296296
70	Double Road	500.0	3.0	100.0	3	2000.000000

```
df3 = df2[~(df2.total_sqft / df2.bhk < 300)]
```

```
df3.shape
```

```
(12582, 6)
```

```
df3.price_per_sqft.describe()
```

count	12456.000000
mean	6308.582826
std	4168.127339
min	267.820913
25%	4218.528316
50%	5294.127647
75%	6916.666667
max	176478.588235
Name: price_per_sqft, dtype: float64	

CREATE A FUN WHICH REMOVES THE OUTLIERS FROM PRICE_PER_SQFT USING STANDARD DEVIATION TECHNIQUE

```
def remove_outlier_from_price_per_sqft(df):
    df_out = pd.DataFrame()
    for key,sub in df.groupby('location'):
        n = np.mean(sub.price_per_sqft)
        st = np.std(sub.price_per_sqft)
        reduce_df = sub[(sub.price_per_sqft>(n-st)) & (sub.price_per_sqft<=(n+st))]
        df_out = pd.concat([df_out, reduce_df],ignore_index=True)
    return df_out
```

```
df4 = remove_outlier_from_price_per_sqft(df3)
```

```
df4.shape
```

```
(9267, 6)
```

```
df4.describe()
```

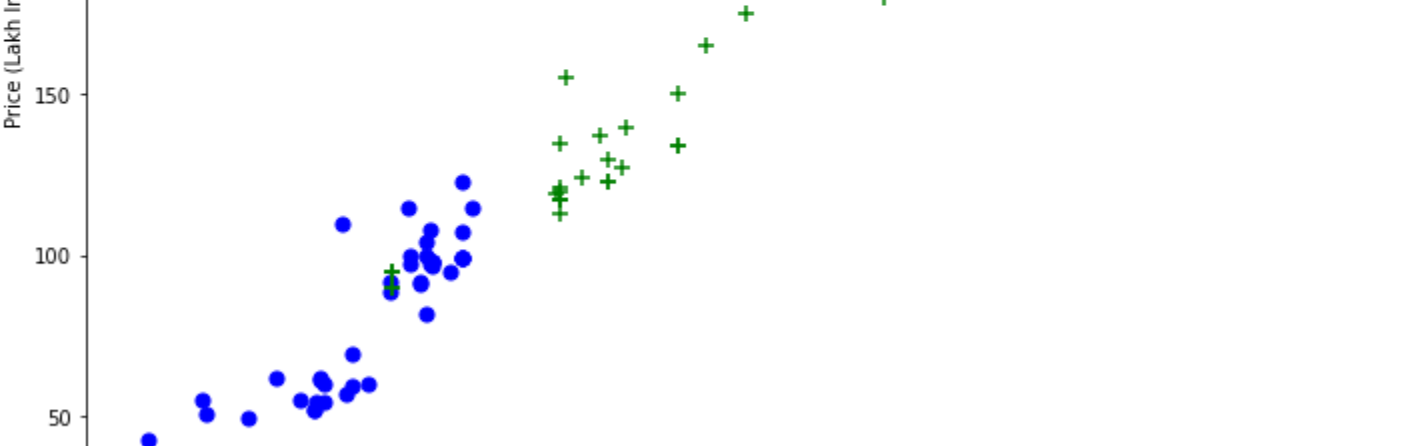
	total_sqft	bath	price	bhk	price_per_sqft
count	9267.000000	9267.000000	9267.000000	9267.000000	9267.000000
mean	2504.002073	2.463149	34.140067	2.556707	5724.037494
std	894.614847	0.953020	110.642062	0.847092	2536.271814
min	300.000000	1.000000	10.000000	1.000000	1250.000000
25%	1109.000000	2.000000	48.000000	2.000000	4299.259259
50%	1262.000000	2.000000	67.000000	2.000000	5185.185185
75%	1650.000000	3.000000	100.000000	3.000000	6004.356285
max	30400.000000	14.000000	2812.000000	10.000000	85000.000000



```
def remove_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, {})
            if stats and stats['count']>5:
                exclude_indices = np.append(exclude_indices, bhk_df[bhk_df.price_per_sqft>(stats['mean']+1)].index.values)
    return df.drop(exclude_indices,axis='index')
df5 = remove_outliers(df4)
df5.shape
```

```
(7564, 6)
```

```
plot_scatter_chart(df5,"Hebbal")
```



```
df5.bath.unique()
```

```
array([ 3., 1., 4., 2., 5., 8., 9., 6., 14., 7., 12.])
```

```
df5[df5.bath==6]
```

	location	total_sqft	bath	price	bhk	price_per_sqft
757	BTM 1st Stage	3300.0	14.0	500.0	9	15151.515152
1951	Chikbaranavali	2400.0	7.0	80.0	4	3250.025020
6117	Nagondandra	7000.0	9.0	450.0	4	6428.571429
7431	Saifys Sai Layout	11338.0	9.0	1000.0	6	8819.897689
7514	Thansandra	1806.0	6.0	116.0	3	6423.034230

```
df6 = df5[(df5.bath > df5.bhk+2)]
```

```
df6.head()
```

	location	total_sqft	bath	price	bhk	price_per_sqft
0	1st Block BEL Layout	1540.0	3.0	85.0	3	5519.480519
1	1st Block HBR Layout	600.0	1.0	45.0	1	7500.000000
2	1st Block HBR Layout	3150.0	4.0	150.0	4	4761.904762
3	1st Block HRBR Layout	2300.0	3.0	80.0	3	3476.260870
4	1st Block HRBR Layout	1250.0	2.0	67.0	2	5360.000000

```
df6.shape
```

```
(7499, 6)
```

```
df7 = df6.drop(['price_per_sqft'],axis='columns')
```

```
df7.head()
```

	location	total_sqft	bath	price	bhk
0	1st Block BEL Layout	1540.0	3.0	85.0	3
1	1st Block HBR Layout	600.0	1.0	45.0	1
2	1st Block HBR Layout	3150.0	4.0	150.0	4
3	1st Block HRBR Layout	2300.0	3.0	80.0	3
4	1st Block HRBR Layout	1250.0	2.0	67.0	2

```
dummies = pd.get_dummies(df7.location)
```

```
dummies.head()
```

	1st Block BEL Layout	1st Block HBR Layout	1st Block HRBR Layout	1st Block Jayanagar	1st Block Koramangala	1st Phase JP Nagar	1st Stage Indira Nagar	2nd Block Jayanagar	2nd Block Jayanagar	Phase JP Nagar	...	Yelahanka New Town	Yelahanka New Town	Yelahanka New Town	Yelenahalli	Yemlur	Yeshwanthpur	Yeshwanthpur Industrial Suburb	frazertown	manyata park	south	tc.palya
0	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0
2	0	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0
3	0	0	1	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0
4	0	0	1	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0

5 rows × 761 columns

```
df8 = pd.concat([df7,dummies.drop('south',axis='columns')],axis='columns')
```

```
df8.head()
```

	location	total_sqft	bath	price	bhk	1st Block BEL Layout	1st Block HBR Layout	1st Block HRBR Layout	1st Block Jayanagar	1st Block Koramangala	...	Yelachenahalli	Yelahanka	Yelahanka New Town	Yelenahalli	Yemlur	Yeshwanthpur	Yeshwanthpur Industrial Suburb	frazertown	manyata park	tc.palya
0	1st Block BEL Layout	1540.0	3.0	85.0	3	1	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	1st Block HBR Layout	600.0	1.0	45.0	1	0	1	0	0	0	...	0	0	0	0	0	0	0	0	0	0