

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression

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

Out[7]:

	area_type	availability	location	size	total_sqft	bath	balcony	price		
0	Super built-up Area	15-Dec	Electronic City Phase II	2 BHK	Comme	1056	2.0	1.0	39.07	
1	Pvt Area	Ready To Move	Chikka Tinsaphi	4 Bedroom	Utarnahe	Thannep	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	Lingdeerenahalli	Ready To Move	Lingdeerenahalli	3 BHK	Sowave	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 [72]:

```
df.shape
```

Out[72]:

```
(13326, 9)
```

In [72]:

```
df.info() #['area_type','availability','society','balcony'],'axis='columns')
df1.head()
```

Out[72]:

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Tinsaphi	4 Bedroom	2600	5.0	120.00
2	Uttarahalli	3 BHK	1440	2.0	62.00
3	Lingdeerenahalli	3 BHK	1521	3.0	95.00
4	Kothanur	2 BHK	1200	2.0	51.00

In [72]:

```
df1.isnull().sum()
```

Out[72]:

```
location      1
total_sqft    0
size          0
bath         15
price        75
dtype: object
```

In [74]:

```
df2=df1.dropna()
df2.isnull().sum()
```

Out[74]:

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

In [75]:

```
df2['size'].unique()
```

Out[75]:

```
array(['2 BHK', '4 Bedroom', '3 BHK', '4 BHK', '6 Bedroom', '3 Bedroom',
       '2 BHK', '1 BHK', '1 Bedroom', '2 Bedroom', '2 Bedroom',
       '7 Bedroom', '5 BHK', '7 BHK', '6 BHK', '5 Bedroom', '11 BHK',
       '9 BHK', '10 BHK', '12 BHK', '10 Bedroom', '12 Bedroom',
       '10 BHK', '13 BHK', '10 BHK', '43 Bedroom', '14 BHK', '8 BHK',
       '12 Bedroom', '13 BHK', '10 Bedroom', dtype=object])
```

In [76]:

```
df2[bhk]>df2['size'].apply(lambda x: int(x.split(' ')[0]))
df2.head()
```

Out[76]:

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
1	Chikka Tinsaphi	4 Bedroom	2600	5.0	120.00	4
2	Uttarahalli	3 BHK	1440	2.0	62.00	3
3	Lingdeerenahalli	3 BHK	1521	3.0	95.00	3
4	Kothanur	2 BHK	1200	2.0	51.00	2

In [77]:

```
df2.total_sqft.unique()
```

Out[77]:

```
array(['1656', '12000', '1440', ..., '1133', '1334', '774', '4689'],
      dtype=object)
```

In [78]:

```
df2.total_sqft
```

Out[78]:

```
0    1056
1     2600
2     1440
3     1521
4     1200
Name: total_sqft, Length: 13246, dtype: object
```

In [79]:

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

Out[79]:

```
df2[df2['total_sqft'].apply(is_float)].head(100)
```

Out[79]:

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
1	Chikka Tinsaphi	4 Bedroom	2600	5.0	120.00	4
2	Uttarahalli	3 BHK	1440	2.0	62.00	3
3	Lingdeerenahalli	3 BHK	1521	3.0	95.00	3
4	Kothanur	2 BHK	1200	2.0	51.00	2

In [80]:

```
df2[df2['total_sqft']>2500]
```

Out[80]:

	location	size	total_sqft	bath	price	bhk
98	Deshnarnagali	2 BHK	1200	2.0	90.00	2
99	T Desaihalli	3 Bedroom	1200	3.0	90.00	3
100	Yeshwanpur	3 BHK	2502	3.0	138.00	3
101	Chandipur	2 BHK	650	1.0	17.00	2
102	Kothanur	3 Bedroom	2400	2.0	150.00	3

100 rows × 6 columns

In [81]:

```
df2.loc[30]
```

Out[81]:

	location	size	total_sqft	bath	price	bhk
30	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
31	Chikka Tinsaphi	4 Bedroom	2600	5.0	120.00	4
32	Uttarahalli	3 BHK	1440	2.0	62.00	3
33	Lingdeerenahalli	3 BHK	1521	3.0	95.00	3
34	Kothanur	2 BHK	1200	2.0	51.00	2

In [82]:

```
df3=df2.copy()
df3['total_sqft']=df3['total_sqft'].apply(pdftotoken)
df3.head()
```

Out[82]:

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tinsaphi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3
3	Lingdeerenahalli	3 BHK	1521.0	3.0	95.00	3
4	Kothanur	2 BHK	1200.0	2.0	51.00	2

In [83]:

```
df4=df3.copy()
len(df4.bath.unique())
```

Out[83]:

```
1394
```

In [86]:

```
df4['price_per_sqft']=df3['price']/df3['total_sqft']
df4.head()
```

Out[86]:

	location	size	total_sqft	bath	price	bhk	price_per_sqft
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2	3699.810066
1	Chikka Tinsaphi	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	Lingdeerenahalli	3 BHK					