

```
import pandas as pd
from matplotlib import pyplot as plt
import numpy as np
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"]=(20,16)
```

In [2]:

```
df1=pd.read_csv('Bengaluru_House_Data.csv')
df1.head()
```

Out[2]:

	area_type	availability	location	size	society	total_sqft	bath	bhkny	price
0	Super built-up Area	19-Dec	Electronic City Phase II	2 BHK	Comeee	1056	2.0	1.0	39.07
1	Plet Area	Ready To Move	Chikka Tinspathi	4 Bedroom	Theamp	2600	5.0	3.0	120.00
2	Built-up Area	Ready To Move	Utaruahalli	3 BHK	N/A	1440	2.0	3.0	62.00
3	Super built-up Area	Ready To Move	Lingadeenarahalli	3 BHK	Salewre	1521	3.0	1.0	95.00
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	N/A	1200	2.0	1.0	51.00

In [3]:

```
df1.shape
```

Out[3]:

```
(13326, 9)
```

In [4]:

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

Out[4]:

```
area_type
Built-up Area    2418
Carpet Area      87
Plot Area        2025
Super built-up Area   8794
Name: area_type, dtype: int64
```

In [5]:

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

Out[5]:

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

In [6]:

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

Out[6]:

```
location      1
size          0
total_sqft    0
bath          0
price         0
dtype: int64
```

In [7]:

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

Out[7]:

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

In [8]:

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

Out[8]:

```
array(['2 BHK', '4 Bedroom', '3 BHK', '4 BHK', '6 Bedroom', '3 Bedrooms',
       '1 BHK', '1 RK', '1 Bedroom', '8 Bedrooe', '2 Bedroom',
       '7 Bedroom', '5 BHK', '7 BHK', '6 BHK', '5 Bedroom', '11 BHK',
       '9 BHK', '9 Bedroom', '22 BHK', '10 Bedroom', '11 Bedroom',
       '19 BHK', '19 BHK', '16 BHK', '43 Bedroom', '14 BHK', '9 BHK',
       '12 Bedroom', '13 BHK', '18 Bedroom'], dtype=object)
```

In [9]:

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

C:\Users\VAIDISH-1\AppData\LocalTemp\ipykernel_17240\222266237.py:1: SettingWithCopyWarning:
A value is being 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
df3[['bhk']] = df3['size'].apply(lambda x: int(x.split(' ')[0]))

In [10]:

```
df3.head()
```

Out[10]:

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

In [11]:

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

Out[11]:

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

In [12]:

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

Out[12]:

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

In [13]:

```
df3.total_sqft.unique()
```

Out[13]:

```
array(['1056', '2089', '1440', ..., '1133', '1384', '774', '4689'],
      dtype=object)
```

In [14]:

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

In [15]:

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

Out[15]:

	location	size	total_sqft	bath	price	bhk
30	Yelenahalli	4 BHK	2100-2900	4.0	195.000	4
122	Hebbal	4 BHK	3067-8156	4.0	477.000	4
137	Sh Phase JP Nagar	2 BHK	1042-1105	2.0	54.005	2
165	Sarjapur	2 BHK	1145-1340	2.0	43.400	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
649	Hebbur Road	2 BHK	1139-1440	2.0	63.770	2
648	Avenue 8 Bedroom	-	412SqMch	8.0	265.000	8
661	Marina	2 BHK	1120-1505	2.0	48.130	2
672	Betevahooli	4 Bedroom	3090-5042	4.0	445.000	4

In [16]:

```
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 [17]:

```
convert_sqft_to_num('2166')
```

Out[17]:

```
2166.0
```

In [18]:

```
convert_sqft_to_num('3067 - 8156')
```

Out[18]:

```
5611.5
```

In [19]:

```
convert_sqft_to_num('34.46sq_Meter')
```

Out[19]:

```
34.46
```

In [20]:

```
dfr=df3.copy()
dfr['total_sqft']= dfr['total_sqft'].apply(convert_sqft_to_num)
dfr.head(3)
```

Out[20]:

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tinspathi	4 Bedroom	2600.0	5.0	120.00	