

Python 3 (ipykernel)

Not Trusted

- File

New NotebookDropdown

Python 3 (ipykernel)

Open...

Make a Copy...

Save as...

Renaming...

Save and CheckpointCtrl+S

Revert to CheckpointDropdown

Munday, June 27, 2022 10:12 AM

Print Preview

Download asDropdown

As a Jupyter Notebook (.ipynb)

HTML (.html)

LaTeX (.tex)

Markdown (.md)

PDF via Jupyter (.pdf)

PDF via LaTeX (.pdf)

reST (.rst)

Python (.py)

Reveal.js slides (.slides.html)

PDF via HTML (.html)

Deploy as

Trust Notebook

Close and Halt

• Edit

Cut Cells

Copy Cells

Paste Cells AboveShift+V

Paste Cells Below

Paste Cells & Replace

Delete Cells

Undo Delete Cells

Split CellCtrl+Shift+Minus

Merge Cell Above

Merge Cell Below

Move Cell Up

Move Cell Down

Edit Notebook Metadata

Find and Replace

Cut Cell Attachments

Copy Cell Attachments

Paste Cell Attachments

Insert Image

• View

Toggle Header

Toggle Toolbar

Toggle Line NumbersShift+L

Cell Toolbar

None

Edit Metadata

Raw Cell Format

Slideshow

Attachments

Tags

• Insert

Insert Cell Above

Insert Cell Below

• Cell

Run CellsCtrl+Enter

Run Cells and Select BelowShift+Enter

Run Cells and Insert BelowAlt+Enter

Run All

Run All Above

Run All Below

Cell Type

Code

Markdown

Raw NBConvert

• Current Outputs

Toggle

Toggle Scrollingshift-Q

Clear

• All Output

Toggle

Toggle Scrolling

Clear

• Kernel

Interrupt

Restart

Restart & Clear Output

Restart & Run All

Reconnect

Shutdown

Change kernel

Python 3 (ipykernel)

• Widgets

Save Notebook Widget State

Clear Notebook Widget State

Download Widget State

Embed Widgets

• Help

User Interface Tour

Keyboard Shortcuts

Edit Keyboard Shortcuts

Notebook Help

Markdown

Python Reference

IPython Reference

NumPy Reference

SciPy Reference

Matplotlib Reference

NumPy Reference

pandas Reference

About

=

==

====

Run

Code

=

In [7]:

```
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 [8]:
```

```
df1 = pd.read_csv(r'C:\Users\jadut3\project-work\Bengaluru_House_Data.csv')
```

Out[8]:

	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	Utarahalli	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 [9]:

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

Out[9]:

```
area_type
Built-up Area      2418
Carpet Area         87
Plot Area          2925
Super built-up Area  8790
Name: area_type, dtype: int64
In [10]:
```

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

Out[10]:

	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	Utarahalli	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

In [25]:

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

Out[25]:

```
location      1
size          16
total_sqft    0
bath          73
price         0
dtype: int64
In [11]:
```

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

Out[11]:

```
location      0
size          0
total_sqft    0
bath          0
price         0
dtype: int64
In [27]:
```

df3.shape

Out[27]:

```
(13246, 5)
```

In [28]:

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

Out[28]:

```
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', '12 Bedroom',
       '10 BHK', '19 BHK', '16 BHK', '43 Bedroom', '14 BHK', '8 BHK',
       '12 Bedroom', '13 BHK', '18 Bedroom'], dtype=object)
```

In [14]:

x

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

```
C:\Users\jadut3\AppData\Local\Temp\ipykernel_6836\2222900254.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
```

```
In [32]:df3['bkh'] = df3['size'].apply(lambda x: int(x.split(' ')[0]))
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
```

Out[32]:

	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	Utarahalli	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
5	Whitefield	2 BHK	1170	2.0	38.00	2
6	Old Airport Road	4 BHK	2732	4.0	204.00	4
7	Rajaji Nagar	4 BHK	3300	4.0	600.00	4
8	Marathahalli	3 BHK	1310	3.0	63.25	3
9	Marathi Bazar	6 Bedroom	1020	6.0	370.00	6

In [16]:

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

Out[16]:

```
array([ 2,  4,  3,  6,  1,  8,  7,  5, 11,  9, 27, 10, 19, 16, 43, 14, 12,
```

In [17]:

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

df3[df3.bkh>20]

Out[17]:

	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

In [31]:

xxxxxx

df4.loc[30]

Out[31]:

```
location      Yelahanka
size          4 BHK
total_sqft    2475.0
bath          4.0
price        186.0
bkh           4
Name: 30, dtype: object
In [32]:
```

xxxxxx

df5=df4.copy()

```
df5['price_per_sqft']=df5['price']/100000/df5['total_sqft']
```

Out[32]:

	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	Utarahalli	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

In [36]:

```
len(df5.location.unique())
```

Out[36]:

```
1384
```

In [40]:

```
df5.location = df5.location.apply(lambda x: x.strip())
```

location_stats = df5.groupby('location')['location'].agg('count').sort_values(ascending=False)

location_stats

Out[40]:

```
location
Whitefield      535
Sarjapur Road  392
Electronic City  304
Kanakapura Road 266
Thansandra      236
...
```

1 Giri Nagar

Kanakapura Road,

Kanakapura main Road

Karnataka Shabarinala

Whitefield

Name: location, Length: 1293, dtype: int64

In [41]:

x

len(location_stats[location_stats<=10])

Out[41]:

```
1052
```

In [45]:

x

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

location_stats_less_than_10

Out[45]:

```
location
Bissapur      10
1st Block Koramangala  10
Gunjur Palya  10
Kalkere       10
Sector 1 HSR Layout  10
...
```

1 Giri Nagar

Kanakapura Road,

Kanakapura main Road

Karnataka Shabarinala

Whitefield

Name: location, Length: 1052, dtype: int64

In [43]:

x

len(df5.location.unique())

Out[43]:

```
1293
```

In [46]:

xxxxxx

```
df5.location = df5.location.apply(lambda x: 'other' if x in location_stats_less_than_10 else x)
```

len(df5.location.unique())

Out[46]:

```
242
```

In [47]:

x

```
df5[df5.total_sqft/df5.bkh<300].head()
```

Out[47]:

	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
50	Murugeshpalya	6 Bedroom	1407.0	4.0	150.0	6	10660.980810
60	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 [48]:

df5.shape

Out[48]:

```
(13246, 7)
```

In [49]:

```
df6=df5[-(df5.total_sqft/df5.bkh<300)]
```

df6.shape

Out[49]:

```
(12562, 7)
```

In [50]:

xxxxxx

```
df6.price_per_sqft.describe()
```

Out[50]:

```
count    12456.000000
mean      6308.502826
std       4188.127339
min       297.829813
25%      4210.526316
50%      5294.176471
75%      6916.666667
max      176478.588235
Name: price_per_sqft, dtype: float64
In [59]:
```

```
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[59]:

```
(10245, 7)
```

In [82]:

```
def plot_scatter_chart(df,location):
```

```
    bhk2 = df[(df.location==location) & (df.bkh==2)]
```

```
    bhk3 = df[(df.location==location) & (df.bkh==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')
```

```
    plt.title(location)
```

```
    plt.legend()
```

```
plot_scatter_chart(df7,"Rajaji Nagar")
```

In [85]:

import matplotlib

```
matplotlib.rcParams["figure.figsize"] = (20,10)
```

```
plt.hist(df7.price_per_sqft,rwidth=0.8)
```

```
plt.xlabel("Price Per Square Feet")
```

```
plt.ylabel("Count")
```

Out[85]:

Text(0, 0.5, 'Count')

In []: