

```
In [1]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import matplotlib
matplotlib.rcParams["figure.figsize"] = (20,10)

In [2]: df1 = pd.read_csv('Bogura_House_Price_Dataset.csv')
df1.head()
```

```
Out[2]:
```

	location	size	sqft	bath	price
0	Namja	2	1056	2	4.69
1	Dhorompur	4	2600	3	1.50
2	Tangra	3	1440	2	10.95
3	Hukmapur	3	1521	2	0.56
4	Boro sorodpur	2	1200	2	5.40

```
In [3]: df1.shape
```

```
Out[3]: (247, 5)
```

```
In [4]: df1.isnull().sum()
```

```
Out[4]: location      0
size              0
sqft              0
bath              0
price             0
dtype: int64
```

```
In [5]: df1.shape
```

```
Out[5]: (247, 5)
```

```
In [6]: df1['size'].unique()
```

```
Out[6]: array([2, 4, 3, 5, 6, 1], dtype=int64)
```

```
In [7]: df1.head()
```

```
Out[7]:
```

	location	size	sqft	bath	price
0	Namja	2	1056	2	4.69
1	Dhorompur	4	2600	3	1.50
2	Tangra	3	1440	2	10.95
3	Hukmapur	3	1521	2	0.56
4	Boro sorodpur	2	1200	2	5.40

```
In [8]: df1['size'].unique()
```

```
Out[8]: array([2, 4, 3, 5, 6, 1], dtype=int64)
```

```
In [9]: df1.sqft.unique()
```

```
Out[9]: array([1056, 2600, 1440, 1521, 1200, 1170, 2732, 3300, 1310, 1020, 1090,
2785, 1080, 1180, 2250, 1175, 1180, 1540, 2770, 680, 1755, 2080,
1767, 510, 1250, 660, 1010, 1151, 1025, 1500, 1407, 840, 4395,
845, 9780, 1160, 3090, 1140, 1220, 1350, 1095, 980, 1398, 1569,
1240, 2089, 1295, 1150, 2511, 480, 4460, 1060, 1320, 1325, 1499,
1665, 789, 1060, 710, 1450, 1296, 2894, 1330, 2302, 650, 2480,
1097, 899, 1030, 1540, 782, 1200, 1415, 1110, 1530, 9700, 2497,
1435, 275, 1427, 2081, 1425, 1470, 1380, 450, 1152, 1550, 480,
785, 770, 1242, 1780, 2144, 1784, 1070, 1846, 1340, 1327, 1186,
1783, 1480, 980, 1285, 912, 1225, 1075, 1252, 1089, 1350, 1297,
1730, 2880, 1595, 1788, 1475, 1580, 1295, 3608, 580, 1415, 1787,
2080, 984, 2480, 1080, 1080, 885, 1153, 1148, 1110, 1290, 1033,
3580, 645, 1644, 1577, 4950, 2420, 880, 1270, 890, 1280, 1188,
3845, 2090, 1162, 1055, 1690, 1464, 700, 1864, 915, 1369, 883,
1684, 2626, 1210, 4111, 1762, 1252, 881, 1420, 1490, 1084, 1015,
1017, 1027, 1069, 1349, 1417, 950, 690, 1053, 1010, 1047, 525,
1850, 1438, 1560, 850, 1113, 1385, 1650], dtype=int64)
```

```
In [10]: df1.loc[31]
```

```
Out[10]: location      Nurul
size              2
sqft             1407
bath              2
price             2.42
price      Name: 31, dtype: object
```

```
In [11]: df5 = df1.copy()
df5['price_per_sqft'] = df5['price']/df5['sqft']
df5.head()
```

```
Out[11]:
```

	location	size	sqft	bath	price	price_per_sqft
0	Namja	2	1056	2	4.69	444.129788
1	Dhorompur	4	2600	3	1.50	57.692308
2	Tangra	3	1440	2	10.95	760.416667
3	Hukmapur	3	1521	2	0.56	36.817883
4	Boro sorodpur	2	1200	2	5.40	450.000000

```
In [12]: len(df5.location.unique())
```

```
Out[12]: 229
```

```
In [13]: len(df5.location)
```

```
Out[13]: 247
```

```
In [14]: df5.location = df5.location.apply(lambda x: x.strip())
location_stats = df5.groupby('location')['location'].agg('count').sort_values(ascending=False)
location_stats
```

```
Out[14]: location      2
Shakpala             2
Barakpur             2
Chandai             2
Bhekra              2
Bhandar Paika       2
Dugaria              1
Doikandi             1
Domonpukur           1
Durgapur             1
Kolgari              1
Name: location, Length: 228, dtype: int64
```

```
In [15]: len(location_stats[location_stats==10])
```

```
Out[15]: 229
```

```
In [16]: location_stats_less_than_10 = location_stats[location_stats<=10]
location_stats_less_than_10
```

```
Out[16]: location      2
Shakpala             2
Barakpur             2
Chandai             2
Bhekra              2
Bhandar Paika       2
Dugaria              1
Doikandi             1
Domonpukur           1
Durgapur             1
Kolgari              1
Name: location, Length: 229, dtype: int64
```

```
In [17]: len(df5.location.unique())
```

```
Out[17]: 229
```

```
In [18]: df5.head(10)
```

```
Out[18]:
```

	location	size	sqft	bath	price	price_per_sqft
0	Namja	2	1056	2	4.69	444.129788
1	Dhorompur	4	2600	3	1.50	57.692308
2	Tangra	3	1440	2	10.95	760.416667
3	Hukmapur	3	1521	2	0.56	36.817883
4	Boro sorodpur	2	1200	2	5.40	450.000000
5	Dhokmohvi	2	1170	2	0.94	80.341880
6	Pultari	4	2732	3	0.50	18.301611
7	Chingapur	5	3300	3	1.45	43.939394
8	Mathura	2	1310	2	0.24	18.320611
9	Bamonpara	2	1020	2	0.20	19.607843

```
In [19]: df5[df5.sqft<df5.size*300].head()
```

```
Out[19]:
```

	location	size	sqft	bath	price	price_per_sqft
0	Namja	2	1056	2	4.69	444.129788
1	Dhorompur	4	2600	3	1.50	57.692308
2	Tangra	3	1440	2	10.95	760.416667
3	Hukmapur	3	1521	2	0.56	36.817883
4	Boro sorodpur	2	1200	2	5.40	450.000000

```
In [20]: df5.shape
```

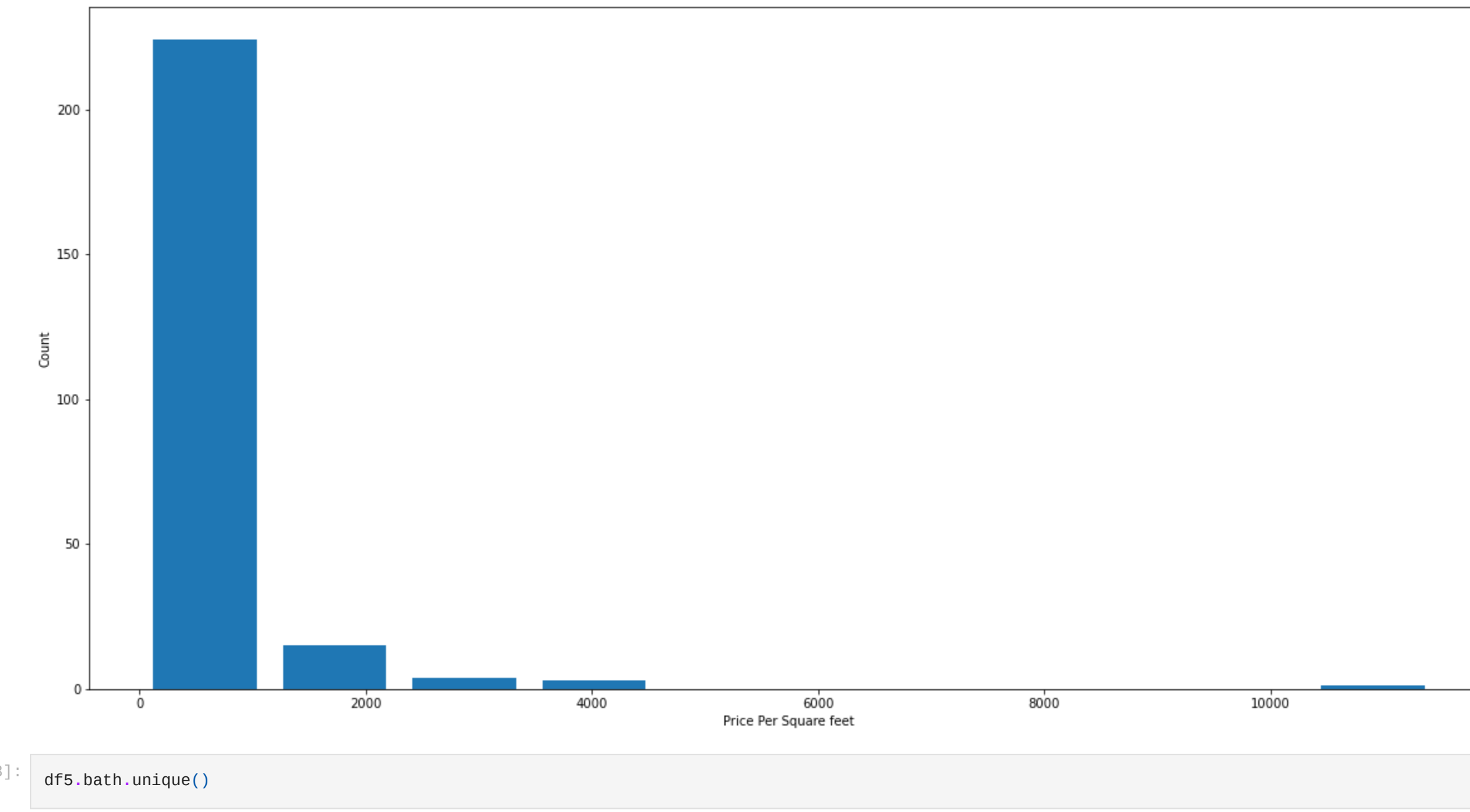
```
Out[20]: (247, 6)
```

```
In [21]: df5.price_per_sqft.describe()
```

```
Out[21]: count      247.000000
mean       387.069747
std        964.911855
min         4.912281
25%        31.899109
50%         90.731707
75%        395.527153
max       11477.079796
Name: price_per_sqft, dtype: float64
```

```
In [22]: import matplotlib
matplotlib.rcParams['figure.figsize'] = (20,10)
plt.hist(df5.price_per_sqft, rwidth=0.8)
plt.xlabel('Price Per Square feet')
plt.ylabel('Count')
```

```
Out[22]: Text(0, 0.5, 'Count')
```



```
In [23]: df5.bath.unique()
```

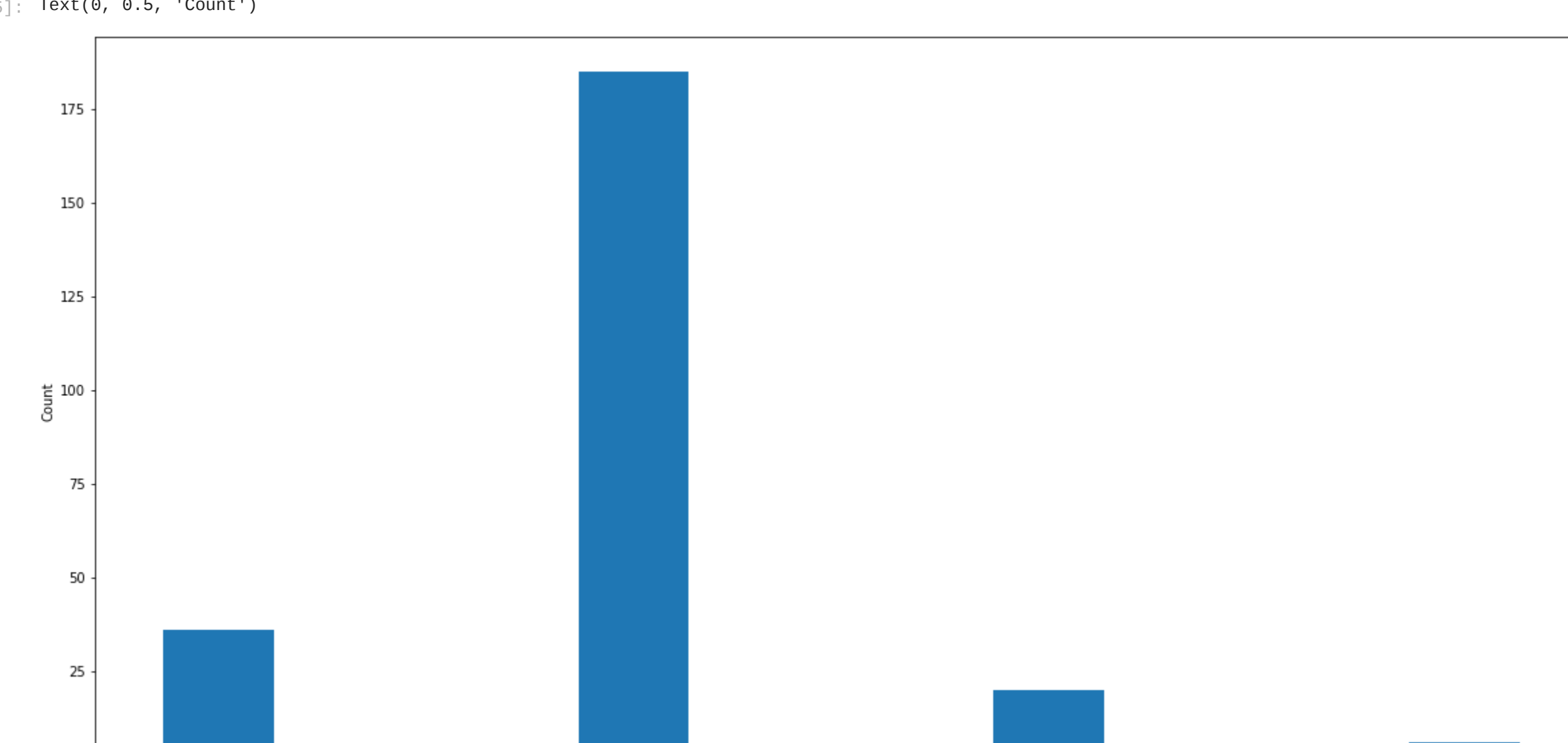
```
Out[23]: array([2, 3, 1, 4], dtype=int64)
```

```
In [24]: # df8(df5.bath=10)
```

```
Out[24]:
```

```
In [25]: plt.hist(df5.bath, rwidth=0.8)
plt.xlabel('Number of bathrooms')
plt.ylabel('Count')
```

```
Out[25]: Text(0, 0.5, 'Count')
```



```
In [26]: # df8(df5.bath=df8.size+2)
```

```
In [27]: # df9 = df8[df8.bath<df8.size+2]
# df9.shape
```

```
Out[28]: df10 = df5.drop(['price_per_sqft'],axis='columns')
df10.head(3)
```

```
Out[28]:
```

	location	size	sqft	bath	price
0	Namja	2	1056	2	4.69
1	Dhorompur	4	2600	3	1.50
2	Tangra	3	1440	2	10.95

```
In [29]: dummies = pd.get_dummies(df10.location)
dummies
```

```
Out[29]:
```

	Agmaji	Agra	Akashtara	Antahar	Aria	Asekpur	Ashokola	Baghopara	Bamonpara	Bamunia	...	Suzabat	Tangra	Tarol	Teldhap	Telhara	Thanthania	Thengamara	Tikshur	Ulipur	Kolgari
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	1	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
...
242	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
243	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
244	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
245	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
246	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

247 rows x 229 columns

```
In [30]: df11 = pd.concat([df10,dummies],axis='columns')
df11.head()
```

```
Out[30]:
```

	location	size	sqft	bath	price	Agmaji	Agra	Akashtara	Antahar	Aria	...	Suzabat	Tangra	Tarol	Teldhap	Telhara	Thanthania	Thengamara	Tikshur	Ulipur	Kolgari
0	Namja	2	1056	2	4.69	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	Dhorompur	4	2600	3	1.50	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
2	Tangra	3	1440	2	10.95	0	0	0	0	0	...	0	1	0	0	0	0	0	0	0	0
3	Hukmapur	3	1521	2	0.56	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
4	Boro sorodpur	2	1200	2	5.40	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

5 rows x 234 columns

```
In [31]: df12 = df11.drop('location',axis='columns')
df12.head(2)
```

```
Out[31]:
```

	size	sqft	bath	price	Agmaji	Agra	Akashtara	Antahar	Aria	Asekpur	...	Suzabat	Tangra	Tarol	Teldhap	Telhara	Thanthania	Thengamara	Tikshur	Ulipur	Kolgari
0	2	1056	2	4.69	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	4	2600	3	1.50	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

2 rows x 233 columns

```
In [32]: df12.shape
# df12
```

```
Out[32]: (247, 233)
```

```
In [33]: X = df12.drop(['price'],axis='columns')
X.head(3)
```

```
Out[33]:
```

	size	sqft	bath	Agmaji	Agra	Akashtara	Antahar	Aria	Asekpur	Ashokola	...	Suzabat	Tangra	Tarol	Teldhap	Telhara	Thanthania	Thengamara	Tikshur	Ulipur	Kolgari
0	2	1056	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	4	2600	3	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
2	3	1440	2	0	0	0	0	0	0	0	...	0	1	0	0	0	0	0	0	0	0

3 rows x 232 columns

```
In [34]: X.shape
X
```

```
Out[34]:
```

	size	sqft	bath	Agmaji	Agra	Akashtara	Antahar	Aria	Asekpur	Ashokola	...	Suzabat	Tangra	Tarol	Teldhap	Telhara	Thanthania	Thengamara	Tikshur	Ulipur	Kolgari
0	2	1056	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	4	2600	3	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
2	3	1440	2	0	0	0	0	0	0	0	...	0	1	0	0	0	0	0	0	0	0
3	3	1521	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
4	2	1200	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
...
242	2	1280	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
243	2	1170	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
244	2	1113	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
245	2	1385	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
246	2	1050	2	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

247 rows x 232 columns

```
In [35]: y = df12.price
y.head(3)
```

```
Out[35]: 0    4.69
1     1.50
2    10.95
Name: price, dtype: float64
```

```
In [36]: 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)
```

```
In [37]: from sklearn.linear_model import LinearRegression
lr_clf = LinearRegression()
lr_clf.fit(X_train,y_train)
lr_clf.score(X_test,y_test)
```

```
Out[37]: 0.815675864173849895
```

```
In [38]: 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)
```

```
Out[38]: array([-5.32838398e-02, -1.19884692e+16, -7.53199276e+17, -8.19076273e+16,
-7.4339738e+02])
```

```
In [39]: def predict_price(location, size, sqft, bath):
    loc_index = np.where(X.columns==location)[0][0]
    # print(loc_index)
    x = np.zeros(len(X.columns))
    x[0] = size
    x[1] = sqft
    x[2] = bath
    if loc_index != 0:
        x[loc_index] = 1
    return round(lr_clf.predict([x])[0], 3)
```

```
In [40]: predict_price('Lotifpur',6,2000,6)
```

```
Out[40]: 51.662
```

```
In [41]: predict_price('Ranirhat', 2, 1050, 2)
```

```
Out[41]: 22.0
```

```
In [42]: predict_price('Ranirhat', 2, 2050, 2)
```

```
Out[42]: 26.894
```

```
In [43]: predict_price('Natal', 2, 1500, 3)
```

```
Out[43]: 36.489
```

```
In [44]: predict_price(location='Ashokola', size=3, sqft=1000, bath=1)
```

```
Out[44]: 0.831
```

```
In [45]: predict_price(location='Lotifpur',size=2,sqft=1000, bath=2)
```

```
Out[45]: 67.056
```

```
In [46]: predict_price('Lotifpur', 2, 1500, 3)
```

```
Out[46]: 66.949
```

```
In [47]: predict_price(location='Ulipur', size=3, sqft=1000, bath=1)
```

```
Out[47]: 0.642
```

```
In [48]: predict_price('Puran Bogra', 3, 1200, 2)
```

```
Out[48]: 11.484
```

```
In [49]: predict_price('Agmaji', 3, 1200, 3)
```

```
Out[49]: 0.053
```

```
In [50]: import pickle
with open('bogura_home_price_model.pickle','wb') as f:
    pickle.dump(lr_clf,f)
```

```
In [51]: import json
columns = {
    'data_columns': [col.lower() for col in X.columns]
}

with open('columns.json','w') as f:
    f.write(json.dumps(columns))
```

```
In [ ]:
```

```
In [ ]:
```