

Import All Libraries

```
In [ ] : import pandas as pd      # pip install pandas
import numpy as np          # pip install numpy
import seaborn as sns       # pip install seaborn
import datetime
import matplotlib.pyplot as plt # pip install matplotlib
import missingno as mso     # pip install missingno
```

Import Dataset

```
In [ ] : df = pd.read_csv('./Data.csv')
```

Dataset Statistics

```
In [ ] : print(f'No of observations (Rows): {df.shape[0]}')
print(f'No of variables (Columns): {df.shape[1]}')
print(f'Duplicate rows: {df.duplicated().sum()}')

nc = 0
for i, j in dict(df.isnull().sum()).items(): nc += int(j) # for getting columns that come with less than 0.5 percent NaN values
print(f'Missing cells: {nc}')
print(f'Missing cells(%): {nc * 100 / (df.shape[0] * df.shape[1])}%')

No of observations (Rows): 191393
No of variables (Columns): 24
Duplicate rows: 0
Missing cells: 94759
Missing cells(%): 2.0629237572255343%
```

Variable Types

```
In [ ] : # Convert columns to best data types
df['date_added'] = pd.to_datetime(df['date_added'])
```

```
In [ ] : print(df.dtypes)
```

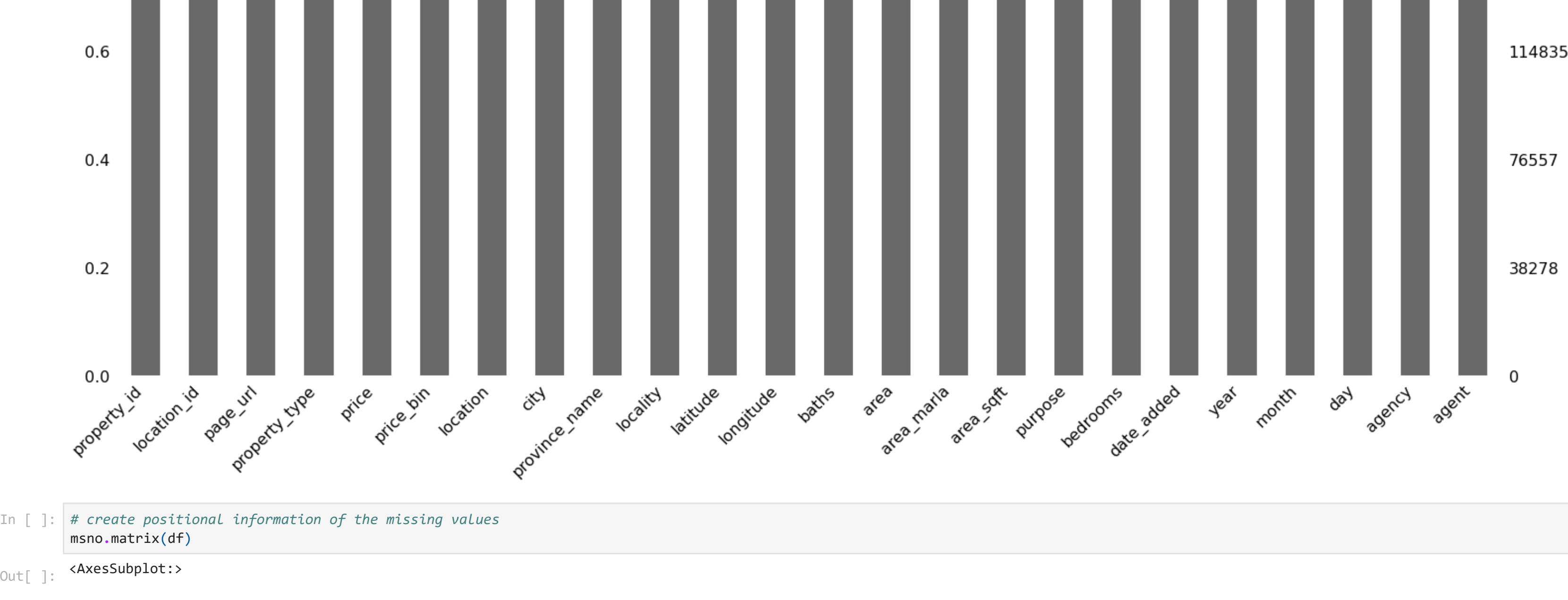
```
property_id      int64
location_id      int64
page_url         object
property_type    object
price            int64
price_bin        object
location         object
city            object
province_name    object
locality         object
latitude         float64
longitude        float64
baths           int64
area            object
area_sqft        object
purpose          object
bedrooms        int64
date_added      datetime64[ns]
year            int64
month           int64
day            int64
agency          object
agent           object
dtype: object
```

```
In [ ] : # count same variables data type
my_dict = {}
list(df.dtypes).count(1) for i in list(df.dtypes)
for i, j in my_dict.items():
    print(f'Variable type: {i}, Count: {j}')
```

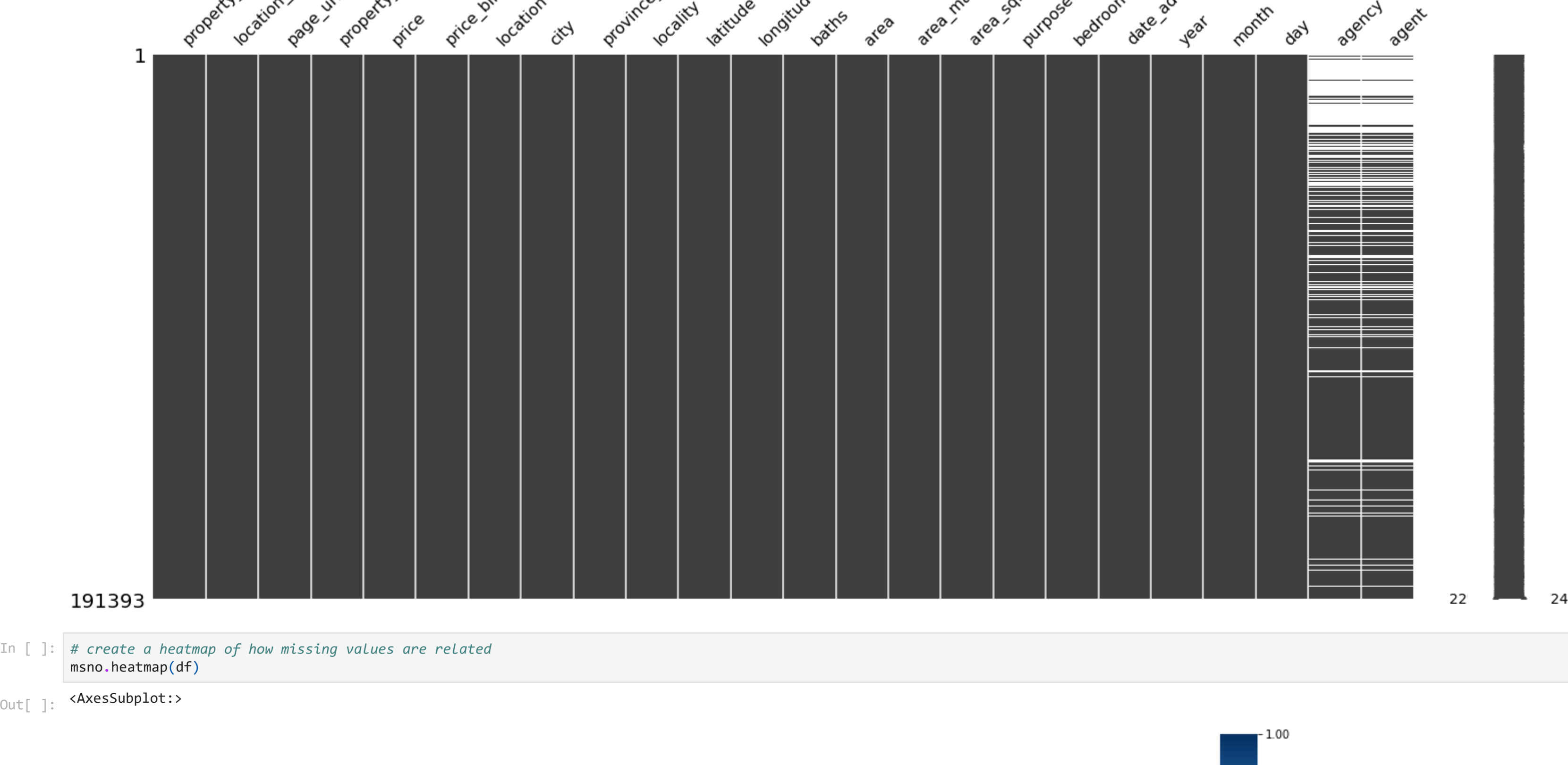
```
Variable type: int64, Count: 8
Variable type: object, Count: 11
Variable type: float64, Count: 4
Variable type: datetime64[ns], Count: 1
```

Missing values

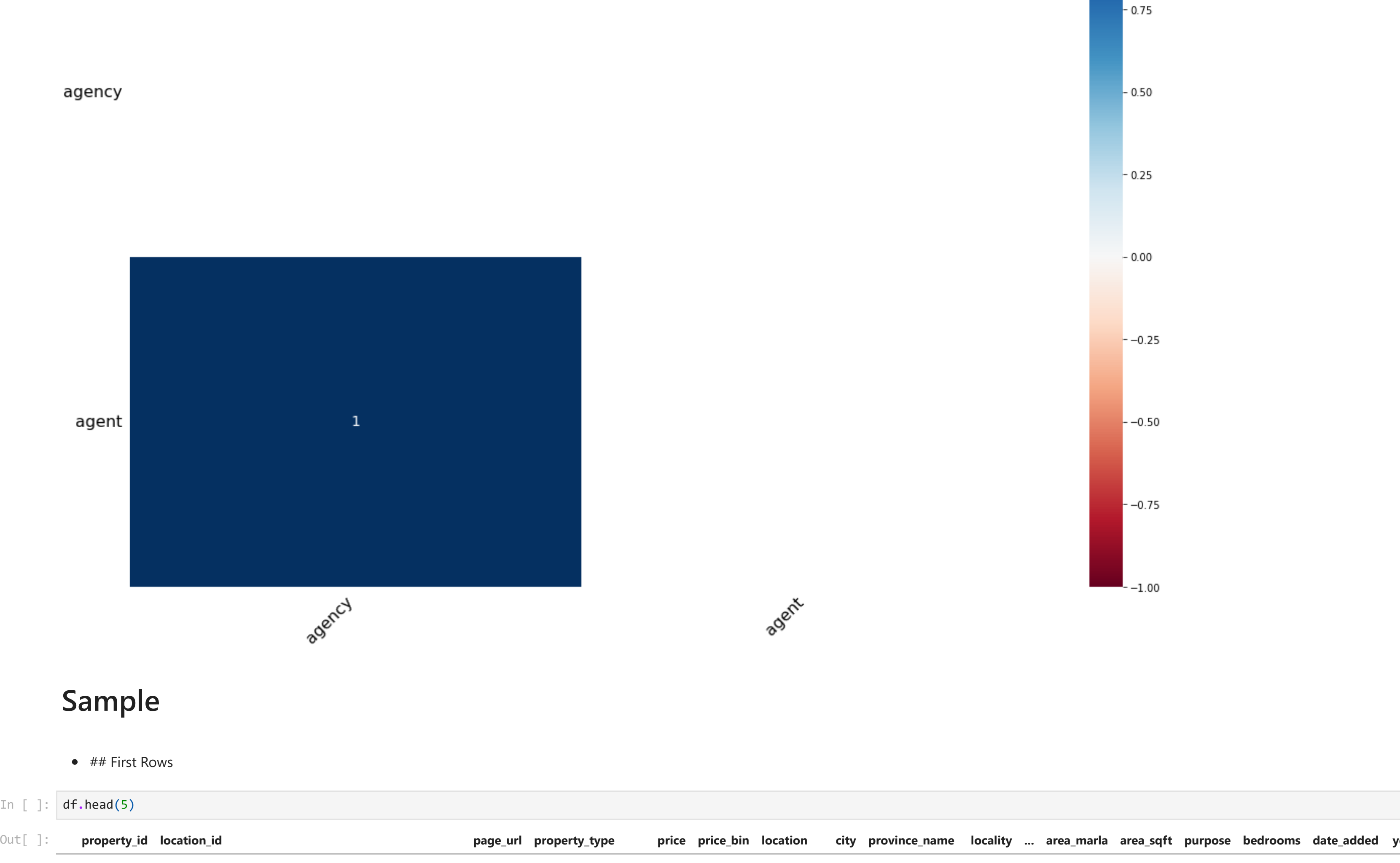
```
In [ ] : # create a bar chart of the missing values
mso.bar(df)
```



```
In [ ] : # create postional information of the missing values
mso.matrix(df)
```



```
In [ ] : # create a heatmap of how missing values are related
mso.heatmap(df)
```



Sample

• ## First Rows

```
In [ ] : df.head(5)
```

	property_id	location_id	page_url	property_type	price	price_bin	location	city	province_name	locality	...	area_marla	area_sqft	purpose	bedrooms	date_added	ye
0	347795	8	https://www.zameen.com/Property/lahore_model_t...	House	220000000	Very High	Model Town	Lahore	Punjab	Model Town, Lahore, Punjab	...	120.0	32670.12	For Sale	0	2019-07-17	20
1	482892	48	https://www.zameen.com/Property/lahore_multan...	House	40000000	Very High	Multan Road	Lahore	Punjab	Multan Road, Lahore, Punjab	...	20.0	5445.02	For Sale	5	2018-10-06	20
2	555962	75	https://www.zameen.com/Property/eden_eden_aven...	House	9500000	Low	Eden	Lahore	Punjab	Eden, Lahore, Punjab	...	9.0	2450.26	For Sale	3	2019-07-03	20
3	562843	3821	https://www.zameen.com/Property/gulberg_2_gulb...	House	125000000	Very High	Gulberg	Lahore	Punjab	Gulberg, Lahore, Punjab	...	20.0	5445.02	For Sale	8	2019-04-04	20
4	686990	3522	https://www.zameen.com/Property/allama_lqbal_t...	House	21000000	High	Allama Iqbal Town	Lahore	Punjab	Allama Iqbal Town, Lahore, Punjab	...	11.0	2994.76	For Sale	6	2019-04-04	20

5 rows × 24 columns

• ## Last Rows

```
In [ ] : df.tail()
```

	property_id	location_id	page_url	property_type	price	price_bin	location	city	province_name	locality	...	area_marla	area_sqft	purpose	bedrooms	date_add
191388	17468383	174	https://www.zameen.com/Property/islamabad_1_8...	Upper Portion	70000	Very High	I-8 Islamabad	Islamabad	Islamabad	I-8, Islamabad, Capital	...	12.4	3375.91	For Rent	3	2019-07-
191389	17468384	174	https://www.zameen.com/Property/islamabad_1_8...	Upper Portion	40000	Medium	I-8 Islamabad	Islamabad	Islamabad	I-8, Islamabad, Capital	...	12.4	3375.91	For Rent	2	2019-07-
191390	17468482	167	https://www.zameen.com/Property/islamabad_g_10...	House	160000	High	G-10 Islamabad	Islamabad	Islamabad	G-10, Islamabad, Capital	...	20.0	5445.02	For Rent	6	2019-07-
191391	17468586	339	https://www.zameen.com/Property/dha_defence_dh...	Flat	25000	Low	DHA Defence	Islamabad	Islamabad	DHA Defence, Islamabad, Capital	...	2.7	735.08	For Rent	2	2019-07-
191392	17468660	3421	https://www.zameen.com/Property/i_10_1_0_2_1...	Upper Portion	26000	Low	I-10 Islamabad	Islamabad	Islamabad	I-10, Islamabad, Capital	...	0.0	0.00	For Rent	3	2019-07-

5 rows × 24 columns

Remove missing values

• ## Remove columns with 95% NaN values

```
In [ ] : for i, j in dict(df.isnull().sum()).items(): # for getting columns that come with less than 0.5 percent NaN values
    if df.shape[0] * j >= 95 * 100 < j: # dropping columns that come with 95% percent NaN values
        df = df.drop(i, axis=1)
```

Out [] :

	property_id	location_id	page_url	property_type	price	price_bin	location	city	province_name	locality	...	area_marla	area_sqft	purpose	bedrooms	date_added	ye
0	347795	8	https://www.zameen.com/Property/lahore_model_t...	House	220000000	Very High	Model Town	Lahore	Punjab	Model Town, Lahore, Punjab	...	120.0	32670.12	For Sale	0	2019-07-17	20
1	482892	48	https://www.zameen.com/Property/lahore_multan...	House	40000000	Very High	Multan Road	Lahore	Punjab	Multan Road, Lahore, Punjab	...	20.0	5445.02	For Sale	5	2018-10-06	20
2	555962	75	https://www.zameen.com/Property/eden_eden_aven...	House	9500000	Low	Eden	Lahore	Punjab	Eden, Lahore, Punjab	...	9.0	2450.26	For Sale	3	2019-07-03	20
3	562843	3821	https://www.zameen.com/Property/gulberg_2_gulb...	House	125000000	Very High	Gulberg	Lahore	Punjab	Gulberg, Lahore, Punjab	...	20.0	5445.02	For Sale	8	2019-04-04	20
4	686990	3522	https://www.zameen.com/Property/allama_lqbal_t...	House	21000000	High	Allama Iqbal Town	Lahore	Punjab	Allama Iqbal Town, Lahore, Punjab	...	11.0	2994.76	For Sale	6	2019-04-04	20

5 rows × 24 columns

• ## Drop NaN values cell

```
In [ ] : df = df.dropna()
df
```

Out [] :

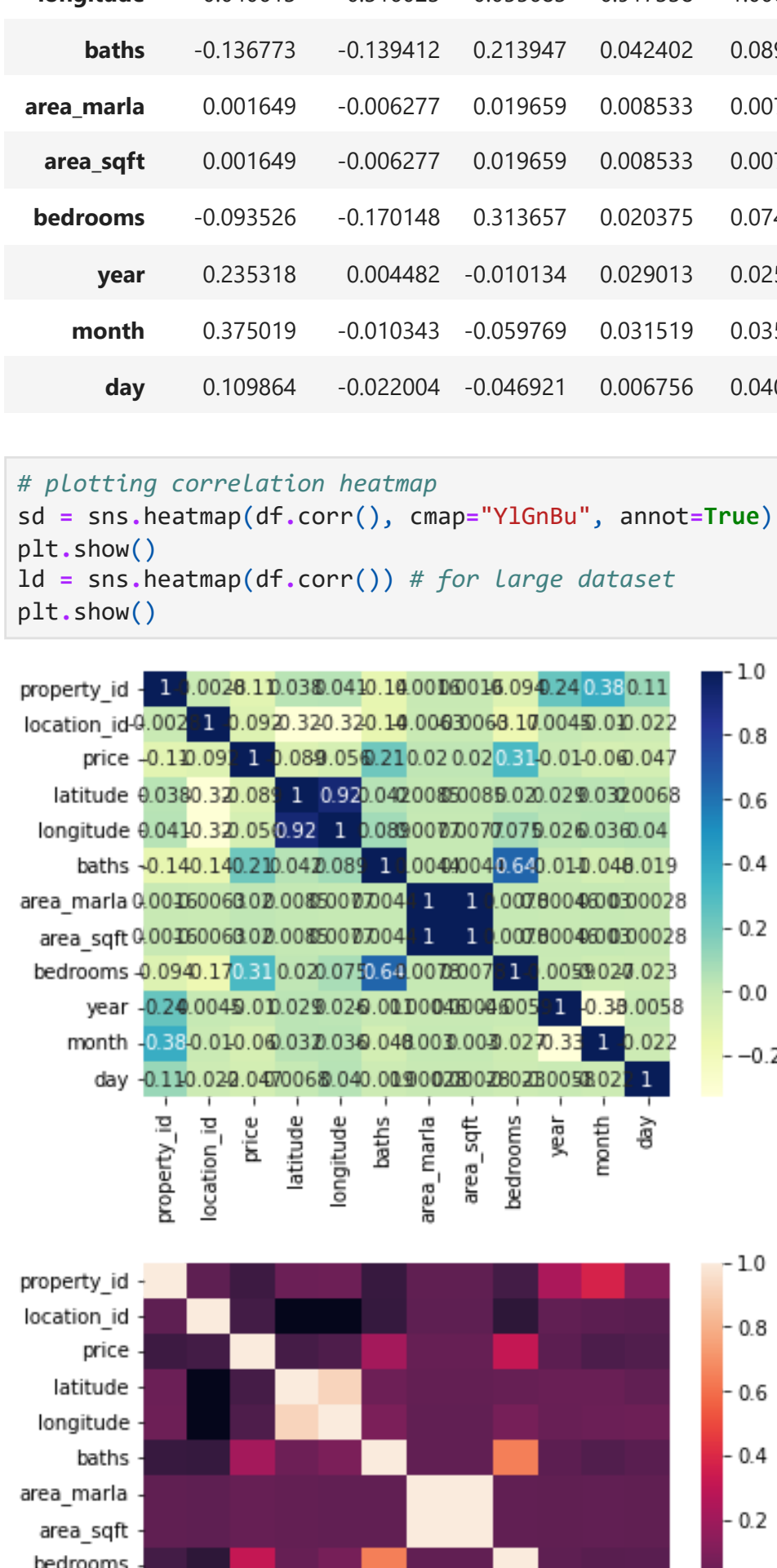
	property_id	location_id	page_url	property_type	price	price_bin	location	city	province_name	locality	...	area_marla	area_sqft	purpose	bedrooms	date
0	347795	8	https://www.zameen.com/Property/lahore_model_t...	House	220000000	Very High	Model Town	Lahore	Punjab	Model Town, Lahore, Punjab	...	120.0	32670.12	For Sale	0	2019
1	482892	48	https://www.zameen.com/Property/lahore_multan...	House	40000000	Very High	Multan Road	Lahore	Punjab	Multan Road, Lahore, Punjab	...	20.0	5445.02	For Sale	5	2018
2	555962	75	https://www.zameen.com/Property/eden_eden_aven...	House	9500000	Low	Eden	Lahore	Punjab	Eden, Lahore, Punjab	...	9.0	2450.26	For Sale	3	2019
5	785289	3102	https://www.zameen.com/Property/gulberg_paf_fa...	House	52000000	Very High	Gulberg	Lahore	Punjab	Gulberg, Lahore, Punjab	...	20.0	5445.02	For Sale	5	2019
7	983055	3749	https://www.zameen.com/Property/eme_society_eme...	House	32500000	High	EME Society	Lahore	Punjab	EME Society, Lahore, Punjab	...	20.0	5445.02	For Sale	5	2019
...
191388	17468383	174	https://www.zameen.com/Property/islamabad_1_8...	Upper Portion	70000	Very High	I-8 Islamabad	Islamabad	Islamabad	I-8, Islamabad, Capital	...	12.4	3375.91	For Rent	3	2019
191389	17468384	174	https://www.zameen.com/Property/islamabad_1_8...	Upper Portion	40000	Medium	I-8 Islamabad	Islamabad	Islamabad	I-8, Islamabad, Capital	...	12.4	3375.91	For Rent	2	2019
191390	17468482	167	https://www.zameen.com/Property/islamabad_g_10...	House	160000	High	G-10 Islamabad	Islamabad	Islamabad	G-10, Islamabad, Capital	...	20.0	5445.02	For Rent	6	2019
191391	17468586	339	https://www.zameen.com/Property/dha_defence_dh...	Flat	25000	Low	DHA Defence	Islamabad	Islamabad	DHA Defence, Islamabad, Capital	...	2.7	735.08	For Rent	2	2019
191392	17468660	3421	https://www.zameen.com/Property/i_10_1_0_2_1...	Upper Portion	26000	Low	I-10 Islamabad	Islamabad	Islamabad	I-10, Islamabad, Capital	...	0.0	0.00	For Rent	3	2019

144013 rows × 24 columns

Correlations

```
In [ ] : df.corr()

sd = sns.heatmap(df.corr(), cmap="YlGnBu", annot=True) # for small dataset
plt.show()
ld = sns.heatmap(df.corr()) # for large dataset
plt.show()
```



Short Details

```
In [ ] : df.info()
```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 144013 entries, 0 to 191392
Data columns (total 24 columns):
Column Non-Null Count Dtype
--
0 property_id 144013 non-null int64
1 location_id 144013 non-null int64
2 page_url 144013 non-null object
3 property_type 144013 non-null object
4 price 144013 non-null int64
5 price_bin 144013 non-null object
6 location 144013 non-null object
7 city 144013 non-null object
8 province_name 144013 non-null object
9 locality 144013 non-null object
10 latitude 144013 non-null float64
11 longitude 144013 non-null float64
12 baths 144013 non-null int64
13 area 144013 non-null object
14 area_marla 144013 non-null float64
15 area_sqft 144013 non-null float64
16 purpose 144013 non-null object
17 bedrooms 144013 non-null int64
18 date_added 144013 non-null datetime64[ns]
19 year 144013 non-null int64
20 month 144013 non-null int64
21 day 144013 non-null int64
22 agency 144013 non-null object
23 agent 144013 non-null object
dtypes: datetime64[ns](1), float64(4), int64(8), object(11)
memory usage: 27.5+ MB

Numerical values short describe

```
In [ ] : df.describe().T # transpose
```

Out [] :

	count	mean	std	min	25%	50%	75%	max
property_id	144013.0	1.634658e+07	1.304215e+06	86575.000000	1.612579e+07	1.696419e+07	1.722522e+07	1.769388e+07
location_id	144013.0	4.186125e+03	6.734307e+03	3.000000	1.244200e+03	3.214000e+03	7.102400e+03	1.424600e+04
price	144013.0	1.756638e+07	3.608480e+07	0.000000	8.500000e+04	7.800000e+06	1.950000e+07	2.000000e+09
latitude	144013.0	3.012832e+01	3.612550e+00	24.749425	2.500530e+01	3.146249e+01	3.355084e+01	7.318409e+01
longitude	144013.0	7.162138e+01	3.065719e+00	66.863657	6.718286e+01	7.307779e+01	7.424677e+01	7.456473e+01
baths	144013.0	2.786748e+00	5.262077e+00	0.000000	0.000000e+00	3.000000e+00	4.000000e+00	4.030000e+02
area_marla	144013.0	1.584714e+01	2.590717e+02	0.000000	5.000000e+00	8.000000e+00	1.400000e+01	1.244440e+05
area_sqft	144013.0	4.314400e+03	1.385953e+05	0.000000	1.361250e+03	2.178010e+03	3.811510e+03	3.388000e+07
bedrooms	144013.0	3.140904e+00	1.932828e+00	0.000000	2.000000e+00	3.000000e+00	4.000000e+00	6.800000e+03
year	144013.0	2.018991e+03	9.360237e-02	0.000000	2.019000e+03	2.019000e+03	2.019000e+03	2.019000e+03
month	144013.0	6.228229e+00	1.238234e+00	1.000000	6.000000e+00	7.000000e+00	7.000000e+00	1.200000e+01
day	144013.0	1.384647e+01	8.675224e+00	1.000000	5.000000e+00	1.400000e+01	2.000000e+01	3.000000e+01