

```
pip install pandas
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (1.5.3)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.10/dist-packages (from pandas) (1.22.4)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2022.7.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)
```

```
pip install seaborn
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (0.12.2)
Requirement already satisfied: matplotlib!=3.6.1,>=3.1 in /usr/local/lib/python3.10/dist-packages (from seaborn) (3.7.1)
Requirement already satisfied: numpy!=1.24.0,>=1.17 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.22.4)
Requirement already satisfied: pandas>=0.25 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.5.3)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.0.7)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (2.8.2)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (0.11.0)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (8.4.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (4.38.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (23.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.4.5)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.25->seaborn) (2022.7.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.1->seaborn) (1.16.0)
```

```
pip install matplotlib
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (3.7.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.4.5)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (2.8.2)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.0.7)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (0.11.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (23.1)
Requirement already satisfied: numpy>=1.20 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.22.4)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (3.0.9)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (4.38.0)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (8.4.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)
```

```
pip install numpy
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (1.22.4)
```

LOAD THE DATASET

```
import pandas as pd
import numpy as np
```

Double-click(or enter) to edit

```
df=pd.read_csv('/content/drive/MyDrive/House Price.csv')
```

```
df
```

	id	Date	number of bedrooms	number of bathrooms	living area	lot area	number of floors	waterfront present	
0	6762810145	42491	5	2.50	3650	9050	2.0	0	
1	6762810635	42491	4	2.50	2920	4000	1.5	0	
2	6762810998	42491	5	2.75	2910	9480	1.5	0	
3	6762812605	42491	4	2.50	3310	42998	2.0	0	

df.dtypes

```
id                int64
Date              int64
number of bedrooms    int64
number of bathrooms  float64
living area          int64
lot area             int64
number of floors      float64
waterfront present    int64
number of views       int64
condition of the house int64
grade of the house    int64
Area of the house(excluding basement) int64
Area of the basement  int64
Built Year            int64
Renovation Year       int64
Postal Code           int64
Lattitude             float64
Longitude             float64
living_area_renov      int64
lot_area_renov         int64
Number of schools nearby int64
Distance from the airport int64
Price                 int64
dtype: object
```

visualization

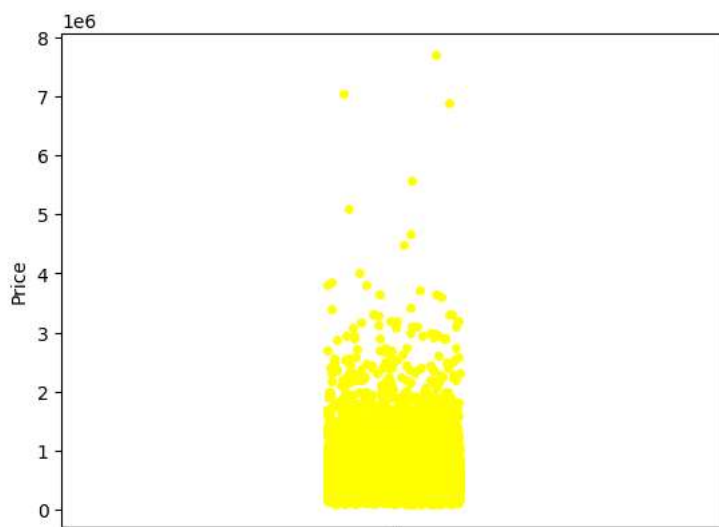
univariate analysis

```
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

```
plt.hist(df['Price'])
```

```
(array([1.2916e+04, 1.4260e+03, 1.9100e+02, 6.0000e+01, 1.9000e+01,
        2.0000e+00, 2.0000e+00, 1.0000e+00, 1.0000e+00, 2.0000e+00]),
array([1.2916e+04, 1.4260e+03, 1.9100e+02, 6.0000e+01, 1.9000e+01,
        2.0000e+00, 2.0000e+00, 1.0000e+00, 1.0000e+00, 2.0000e+00]))
sns.stripplot(y=df['Price'],color='yellow')
```

<Axes: ylabel='Price'>



BI-VARIATE ANALYSIS

```
rate=pd.read_csv('/content/drive/MyDrive/House Price.csv')
rate.plot(x='id',y='price',kind='scatter',color='indigo');
plt.show
```

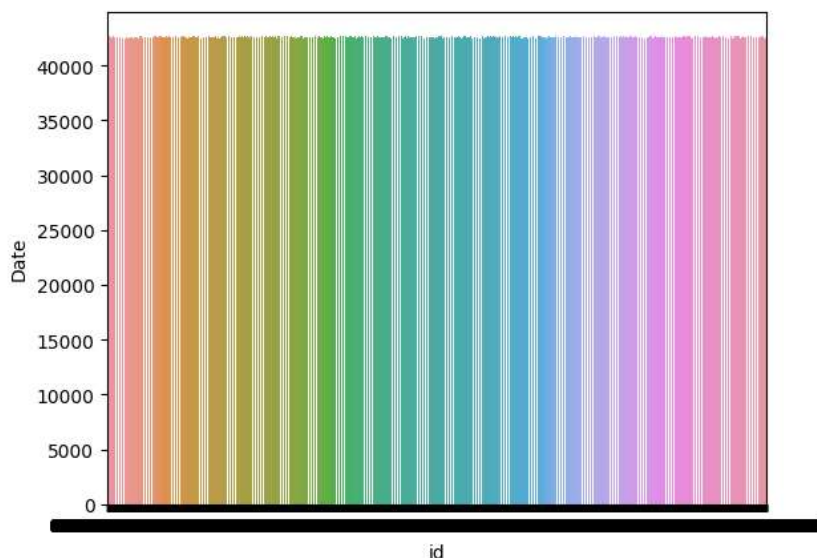
```
-----
NameError                                Traceback (most recent call last)
<ipython-input-2-5d90e80866f4> in <cell line: 1>()
----> 1 rate=pd.read_csv('/content/drive/MyDrive/House Price.csv')
      2 rate.plot(x='id',y='price',kind='scatter',color='indigo');
      3 plt.show
```

NameError: name 'pd' is not defined

SEARCH STACK OVERFLOW

```
sns.barplot(data=rate,x='id',y='Date');
plt.show
```

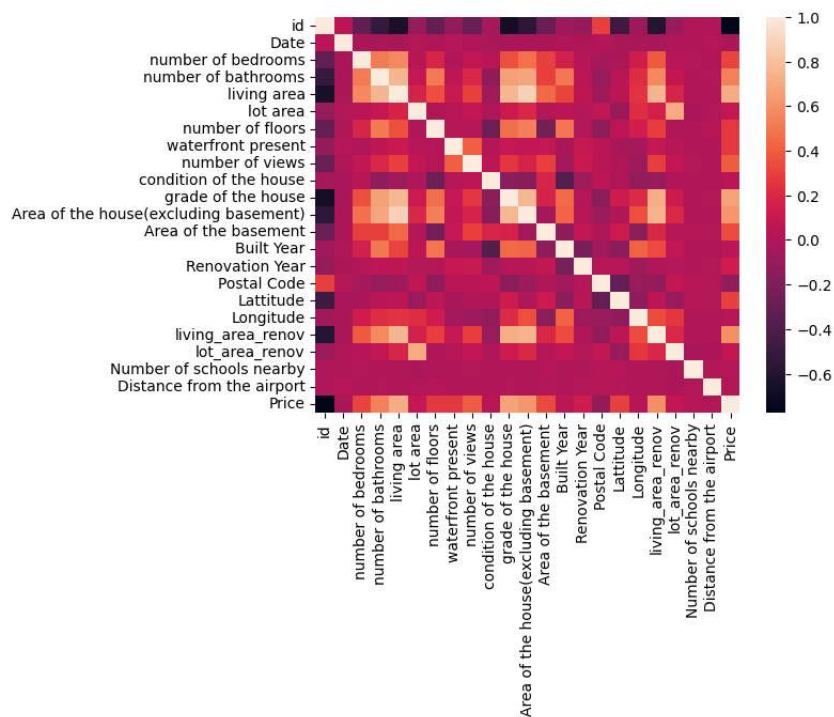
<function matplotlib.pyplot.show(close=None, block=None)>



MULTI-VARIATE ANALYSIS

```
sns.heatmap(rate.corr(),annot=False)
```

<Axes: >



DESCRIPTIVE STATISTICS

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
import warnings
warnings.filterwarnings('ignore')
```

```
data='/content/drive/MyDrive/House Price.csv'
df=pd.read_csv(data)
```

```
df.shape

(14620, 23)
```

```
df.head()
```

```
df.describe(include='all')
```

	id	Date	number of bedrooms	number of bathrooms	living area	lot area	number of floors	waterfront present	number of views
count	1462	1462	1462	1462	1462	1462	1462	1462	1462
mean	6762810145	42491	5	2.50	3650	9050	2.0	0	1.509321
std	6.237575e+03	67.347991	0.938719	0.769934	928.275721	3.791961	0.000000	0.000000	0.000000
min	6.762810e+09	42491.000000	1.000000	0.500000	370.000000	5.200000	0.000000	0.000000	0.000000
25%	6.762815e+09	42546.000000	3.000000	1.750000	1440.000000	5.010751	0.000000	0.000000	0.000000
50%	6.762821e+09	42600.000000	3.000000	2.250000	1930.000000	7.620000	0.000000	0.000000	0.000000
75%	6.762826e+09	42662.000000	4.000000	2.500000	2570.000000	1.080000	0.000000	0.000000	0.000000
max	6.762832e+09	42734.000000	33.000000	8.000000	13540.000000	1.074211	0.000000	0.000000	0.000000

MEAN

```
mean=df['waterfront present'].mean()
print(mean)
```

0.007660738714090287

MEDIAN

```
median=df['waterfront present'].median()
print(median)
```

0.0

MODE

```
mode=df['waterfront present'].mode()
print(mode)
```

0 0
Name: waterfront present, dtype: int64

OBSERVATION

PLOT THE DISTRIBUTION

```
data=df['waterfront present']
sns.distplot(data,bins=10,hist=True,kde=True,label='waterfront present')
```

<Axes: xlabel='waterfront present', ylabel='Density'>



check for missing values



```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14620 entries, 0 to 14619
Data columns (total 23 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     14620 non-null  int64
1   Date                                 14620 non-null  int64
2   number of bedrooms                   14620 non-null  int64
3   number of bathrooms                  14620 non-null  float64
4   living area                          14620 non-null  int64
5   lot area                             14620 non-null  int64
6   number of floors                     14620 non-null  float64
7   waterfront present                   14620 non-null  int64
8   number of views                      14620 non-null  int64
9   condition of the house               14620 non-null  int64
10  grade of the house                   14620 non-null  int64
11  Area of the house(excluding basement) 14620 non-null  int64
12  Area of the basement                 14620 non-null  int64
13  Built Year                           14620 non-null  int64
14  Renovation Year                      14620 non-null  int64
15  Postal Code                          14620 non-null  int64
16  Lattitude                            14620 non-null  float64
17  Longitude                            14620 non-null  float64
18  living_area_renov                    14620 non-null  int64
19  lot_area_renov                       14620 non-null  int64
20  Number of schools nearby              14620 non-null  int64
21  Distance from the airport            14620 non-null  int64
22  Price                                14620 non-null  int64
dtypes: float64(4), int64(19)
memory usage: 2.6 MB
```

```
print(df.isnull().sum())
```

```
id                                     0
Date                                 0
number of bedrooms                   0
number of bathrooms                  0
living area                          0
lot area                             0
number of floors                     0
waterfront present                   0
number of views                      0
condition of the house               0
grade of the house                   0
Area of the house(excluding basement) 0
Area of the basement                 0
Built Year                           0
Renovation Year                      0
Postal Code                          0
Lattitude                            0
Longitude                            0
living_area_renov                    0
lot_area_renov                       0
Number of schools nearby              0
Distance from the airport            0
Price                                0
dtype: int64
```

```
update_df=df.dropna(axis=1)
```

```
update_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14620 entries, 0 to 14619
Data columns (total 23 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   id                                         14620 non-null  int64
1   Date                                       14620 non-null  int64
2   number of bedrooms                       14620 non-null  int64
3   number of bathrooms                     14620 non-null  float64
4   living area                              14620 non-null  int64
5   lot area                                 14620 non-null  int64
6   number of floors                         14620 non-null  float64
7   waterfront present                       14620 non-null  int64
8   number of views                         14620 non-null  int64
9   condition of the house                  14620 non-null  int64
10  grade of the house                      14620 non-null  int64
11  Area of the house(excluding basement)    14620 non-null  int64
12  Area of the basement                    14620 non-null  int64
13  Built Year                              14620 non-null  int64
14  Renovation Year                         14620 non-null  int64
15  Postal Code                             14620 non-null  int64
16  Latitude                                14620 non-null  float64
17  Longitude                               14620 non-null  float64
18  living_area_renov                       14620 non-null  int64
19  lot_area_renov                         14620 non-null  int64
20  Number of schools nearby                 14620 non-null  int64
21  Distance from the airport                14620 non-null  int64
22  Price                                    14620 non-null  int64
dtypes: float64(4), int64(19)
memory usage: 2.6 MB
```

```
update_df=df
```

```
update_df['waterfront present']=update_df['waterfront present'].fillna(update_df['waterfront present'].mean())
```

```
update_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14620 entries, 0 to 14619
Data columns (total 23 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   id                                         14620 non-null  int64
1   Date                                       14620 non-null  int64
2   number of bedrooms                       14620 non-null  int64
3   number of bathrooms                     14620 non-null  float64
4   living area                              14620 non-null  int64
5   lot area                                 14620 non-null  int64
6   number of floors                         14620 non-null  float64
7   waterfront present                       14620 non-null  int64
8   number of views                         14620 non-null  int64
9   condition of the house                  14620 non-null  int64
10  grade of the house                      14620 non-null  int64
11  Area of the house(excluding basement)    14620 non-null  int64
12  Area of the basement                    14620 non-null  int64
13  Built Year                              14620 non-null  int64
14  Renovation Year                         14620 non-null  int64
15  Postal Code                             14620 non-null  int64
16  Latitude                                14620 non-null  float64
17  Longitude                               14620 non-null  float64
18  living_area_renov                       14620 non-null  int64
19  lot_area_renov                         14620 non-null  int64
20  Number of schools nearby                 14620 non-null  int64
21  Distance from the airport                14620 non-null  int64
22  Price                                    14620 non-null  int64
dtypes: float64(4), int64(19)
memory usage: 2.6 MB
```

```
update_df=df
```

```
update_df['waterfront present missing']=update_df['waterfront present'].isnull()
```

```
from sklearn.impute import SimpleImputer
```

```
my_imputer=SimpleImputer(strategy='median')
```

```
data_new=my_imputer.fit_transform(update_df)
```

```
update_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14620 entries, 0 to 14619
Data columns (total 24 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   id                                         14620 non-null  int64
```

```
1 Date 14620 non-null int64
2 number of bedrooms 14620 non-null int64
3 number of bathrooms 14620 non-null float64
4 living area 14620 non-null int64
5 lot area 14620 non-null int64
6 number of floors 14620 non-null float64
7 waterfront present 14620 non-null int64
8 number of views 14620 non-null int64
9 condition of the house 14620 non-null int64
10 grade of the house 14620 non-null int64
11 Area of the house(excluding basement) 14620 non-null int64
12 Area of the basement 14620 non-null int64
13 Built Year 14620 non-null int64
14 Renovation Year 14620 non-null int64
15 Postal Code 14620 non-null int64
16 Latitude 14620 non-null float64
17 Longitude 14620 non-null float64
18 living_area_renov 14620 non-null int64
19 lot_area_renov 14620 non-null int64
20 Number of schools nearby 14620 non-null int64
21 Distance from the airport 14620 non-null int64
22 Price 14620 non-null int64
23 waterfront present missing 14620 non-null bool
dtypes: bool(1), float64(4), int64(19)
memory usage: 2.6 MB
```