

Linear Regression

```
In [1]: import pandas as pd
import numpy as np

In [2]: import matplotlib.pyplot as plt
import seaborn as sns

In [3]: %matplotlib inline

In [5]: df = pd.read_csv('USA_Housing.csv')

In [6]: df.head()

Out[6]:
   Avg. Area Income  Avg. Area House Age  Avg. Area Number of Rooms  Avg. Area Number of Bedrooms  Area Population  Price  Address
0      79545.458574          5.682861             7.009188                4.09      23086.800503  1.059034e+06  208 Michael Ferry Apt. 674inLaurabury, NE 3701...
1      79248.642455          6.002900             6.730821                3.09      40173.072174  1.505891e+06  188 Johnson Views Suite 079nLake Kathleen, CA...
2      61287.067179          5.865890             8.512727                5.13      36882.159400  1.058988e+06  9127 Elizabeth StravenueinDanieltown, WI 06482...
3      63345.240046          7.188236             5.586729                3.26      34310.242831  1.260617e+06                USS BarnettinFPO AP 44820
4      59982.197226          5.040555             7.839388                4.23      26354.109472  6.309435e+05                USNS RaymondinFPO AE 09386

In [7]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5090 entries, 0 to 4999
Data columns (total 7 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Avg. Area Income    5090 non-null  float64
1   Avg. Area House Age 5090 non-null  float64
2   Avg. Area Number of Rooms 5090 non-null float64
3   Avg. Area Number of Bedrooms 5090 non-null float64
4   Area Population      5090 non-null float64
5   Price               5090 non-null float64
6   Address             5090 non-null object
dtypes: float64(6), object(1)
memory usage: 273.6+ KB

In [8]: df.describe()

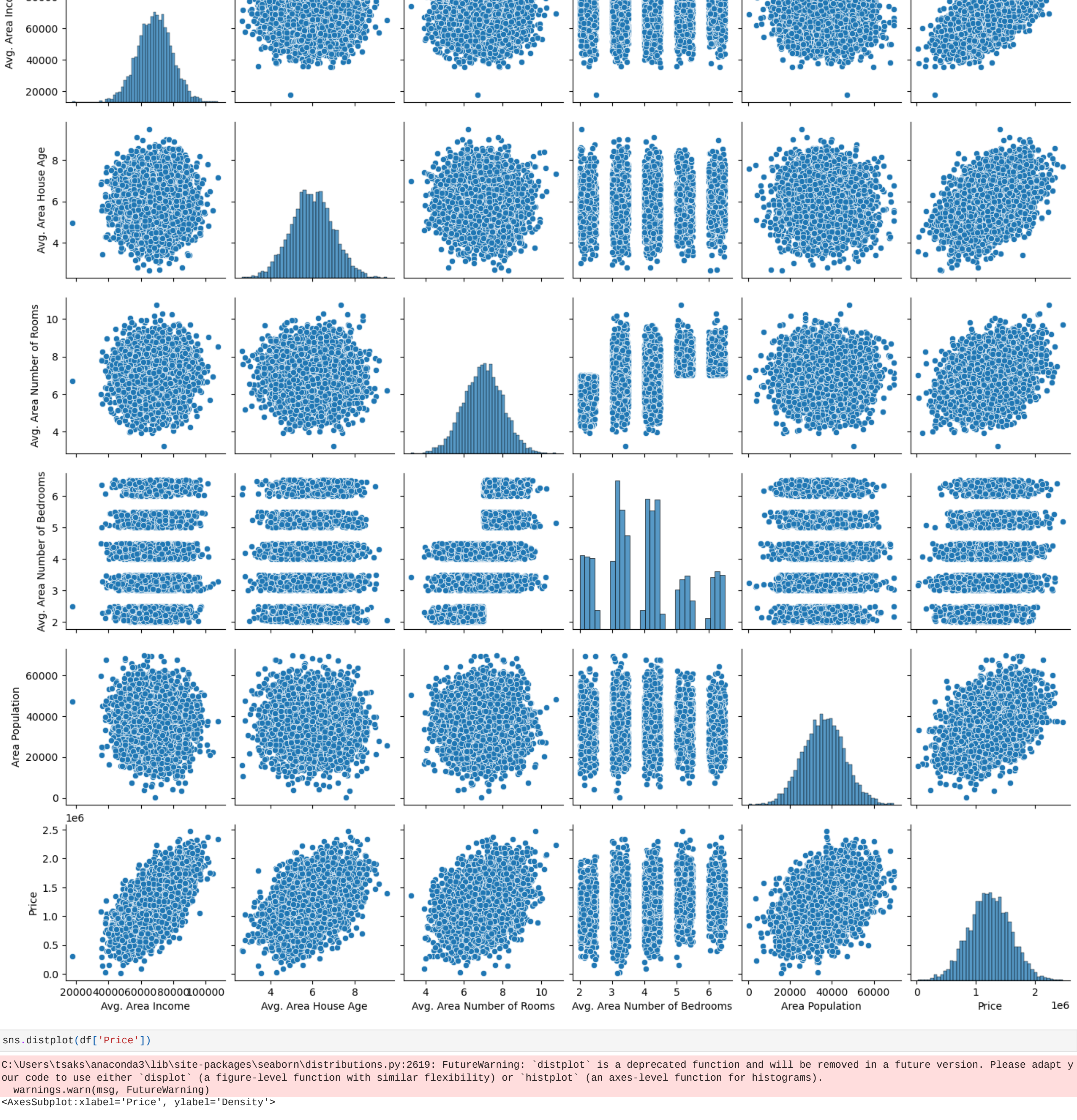
Out[8]:
           Avg. Area Income  Avg. Area House Age  Avg. Area Number of Rooms  Avg. Area Number of Bedrooms  Area Population  Price
count      5000.000000          5000.000000          5000.000000          5000.000000      5000.000000  5.000000e+03
mean      68583.108984          5.977222             6.987792                3.981330      36163.516039  1.232073e+06
std       10657.991214          0.991456             1.005833                1.234137      9925.650114  3.531176e+05
min       17796.631190          2.644304             3.236194                2.000000      172.610686  1.593866e+04
25%       61480.562388          5.322283             6.299250                3.140000      29403.928702  9.975771e+05
50%       68804.286404          5.970429             7.002902                4.050000      36199.406689  1.232669e+06
75%       75783.338666          6.650808             7.665871                4.490000      42861.290769  1.471210e+06
max       107701.748378          9.519088             10.759588                6.500000      69621.713378  2.469066e+06

In [9]: df.columns

Out[9]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
        'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],
        dtype='object')

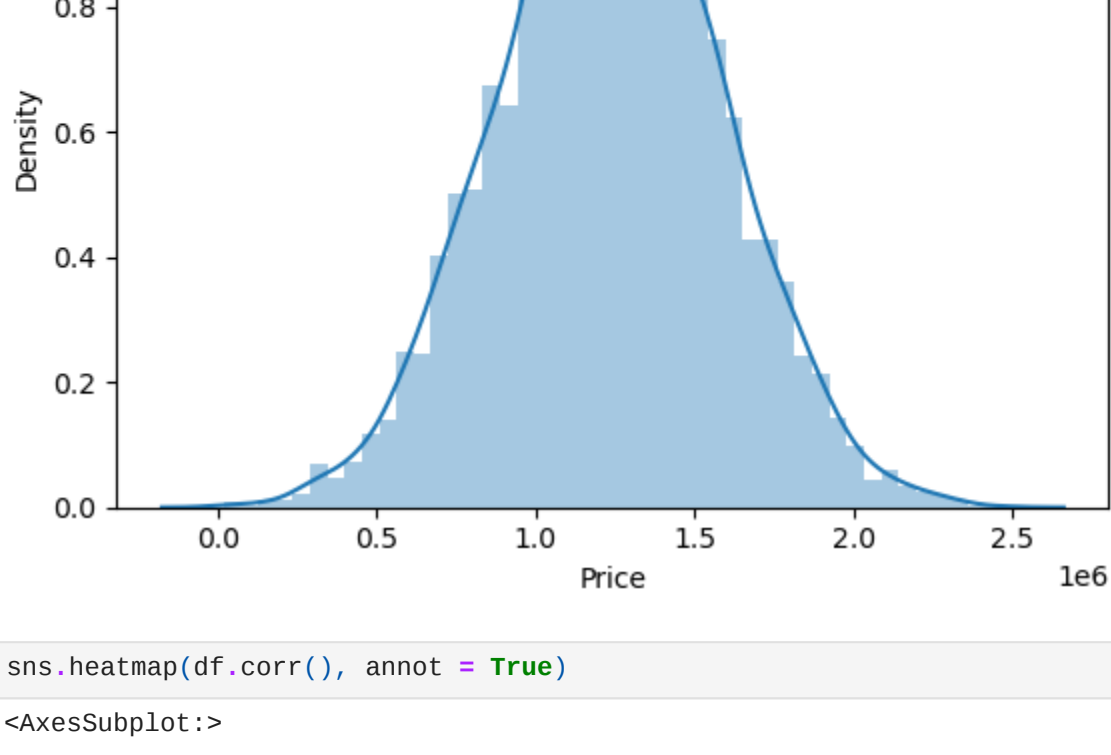
In [10]: sns.pairplot(df)

Out[10]: <seaborn.axisgrid.PairGrid at 0x17870c739a8>
```



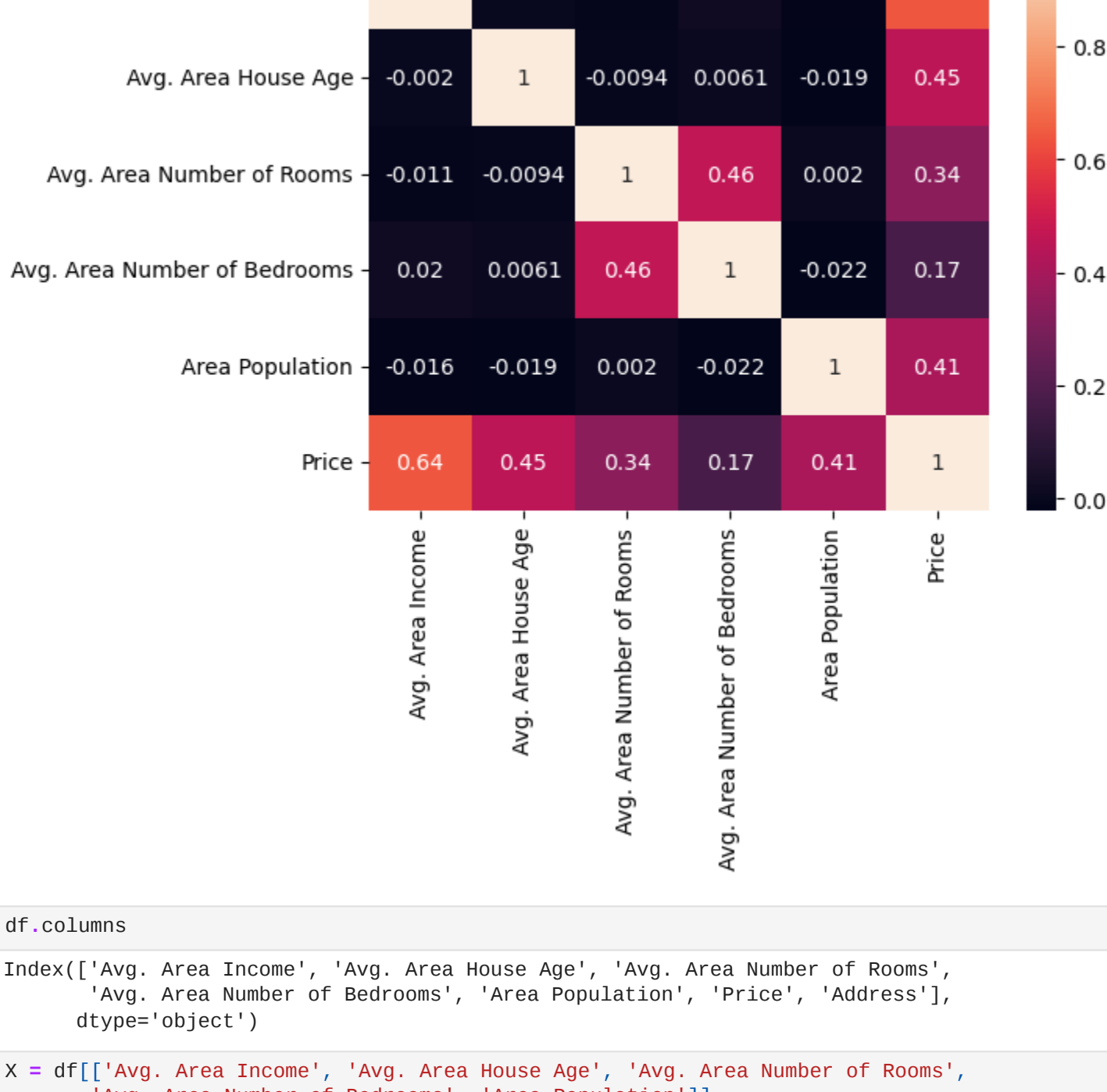
```
In [12]: sns.distplot(df['Price'])

C:\Users\tsaks\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt y
our code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
<AxesSubplot: xlabel='Price', ylabel='Density'>
```



```
In [14]: sns.heatmap(df.corr(), annot = True)

Out[14]: <AxesSubplot: >
```



```
In [15]: df.columns

Out[15]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
        'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],
        dtype='object')

In [16]: X = df[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
        'Avg. Area Number of Bedrooms', 'Area Population']]

In [17]: y = df['Price']

In [18]: from sklearn.model_selection import train_test_split

In [19]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=181)

In [20]: from sklearn.linear_model import LinearRegression

In [21]: lm = LinearRegression()

In [22]: lm.fit(X_train,y_train)

Out[22]: LinearRegression()

In [23]: print(lm.intercept_)

-2648159.7968525267

In [24]: lm.coef_

array([[2.15282755e+01, 1.64883282e+05, 1.22368678e+05, 2.23380186e+03,
        1.51504200e+01]])

In [25]: X_train.columns

Out[25]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
        'Avg. Area Number of Bedrooms', 'Area Population'],
        dtype='object')

In [27]: cdf = pd.DataFrame(lm.coef_,X.columns, columns=['Coeff'])

In [28]: cdf

Out[28]:
           Coeff
Avg. Area Income      21.528276
Avg. Area House Age  164883.282027
Avg. Area Number of Rooms  122368.678027
Avg. Area Number of Bedrooms  2233.801864
Area Population      15.150420
```

Predictions

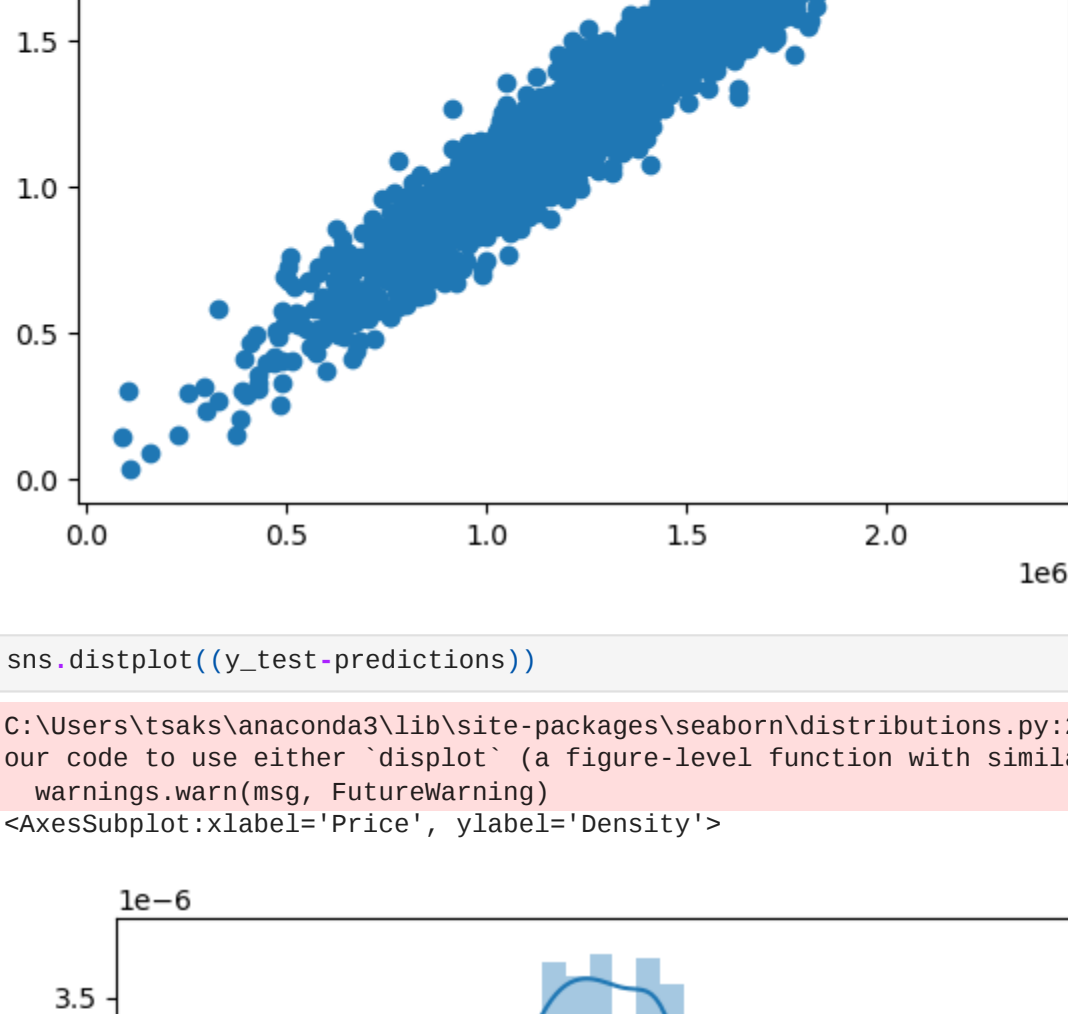
```
In [29]: predictions = lm.predict(X_test)

In [30]: predictions

Out[30]: array([[1260960.70567627, 827588.75560334, 1742421.2425434, ...,
        372191.48626923, 1365217.15140897, 1914519.5417887 ]])

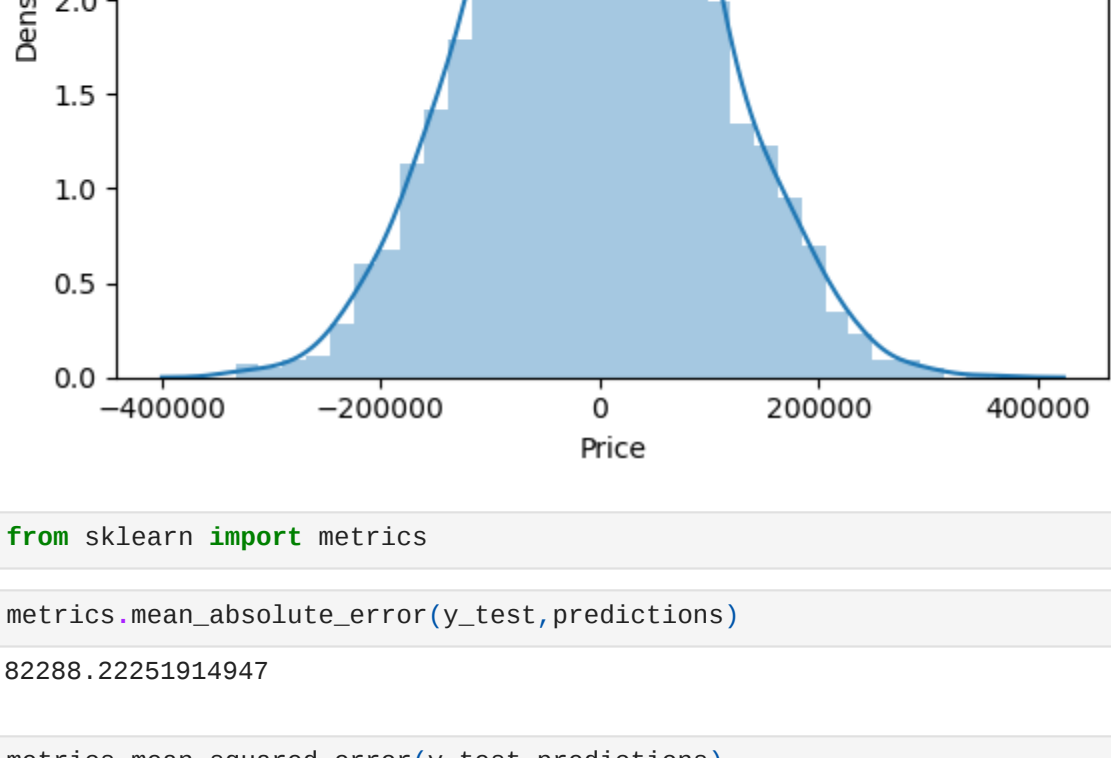
In [31]: plt.scatter(predictions,y_test)

Out[31]: <matplotlib.collections.PathCollection at 0x17876ef370>
```



```
In [32]: sns.distplot(y_test-predictions)

C:\Users\tsaks\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt y
our code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
<AxesSubplot: xlabel='Price', ylabel='Density'>
```



```
In [33]: from sklearn import metrics

In [34]: metrics.mean_absolute_error(y_test,predictions)

Out[34]: 82288.22251914947

In [35]: metrics.mean_squared_error(y_test,predictions)

Out[35]: 10466958907.209059

In [36]: np.sqrt(metrics.mean_squared_error(y_test,predictions))

Out[36]: 102278.82922290936

In [ ]:
```