

Problem statement

A real estate agent want help to predict the house price for regions in USA. He gave us the dataset to work on to use Linear Regression model. Create a model that helps to determine

Data collection

Importing libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Importing dataset

```
In [2]: data=pd.read_csv(r"C:\Users\user\Downloads\H.csv")
data
```

Out[2]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.45857	5.682861	7.009188	4.09	23086.80050	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.64245	6.002900	6.730821	3.09	40173.07217	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.06718	5.865890	8.512727	5.13	36882.15940	1.058988e+06	9127 Elizabeth Stravenue\nDanieltown, WI 06482...
3	63345.24005	7.188236	5.586729	3.26	34310.24283	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.19723	5.040555	7.839388	4.23	26354.10947	6.309435e+05	USNS Raymond\nFPO AE 09386
...
4995	60567.94414	7.830362	6.137356	3.46	22837.36103	1.060194e+06	USNS Williams\nFPO AP 30153-7653
4996	78491.27543	6.999135	6.576763	4.02	25616.11549	1.482618e+06	PSC 9258, Box 8489\nAPO AA 42991- 3352

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
4997	63390.68689	7.250591	4.805081	2.13	33266.14549	1.030730e+06	4215 Tracy Garden Suite 076\nJoshualand, VA 01...
4998	68001.33124	5.534388	7.130144	5.44	42625.62016	1.198657e+06	USS Wallace\nFPO AE 73316
4999	65510.58180	5.992305	6.792336	4.07	46501.28380	1.298950e+06	37778 George Ridges Apt. 509\nEast Holly, NV 2...

5000 rows × 7 columns

head

```
In [3]: # to display first 8 dataset values
data.head(8)
```

Out[3]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.45857	5.682861	7.009188	4.09	23086.80050	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.64245	6.002900	6.730821	3.09	40173.07217	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.06718	5.865890	8.512727	5.13	36882.15940	1.058988e+06	9127 Elizabeth Stravenue\nDanieltown, WI 06482...
3	63345.24005	7.188236	5.586729	3.26	34310.24283	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.19723	5.040555	7.839388	4.23	26354.10947	6.309435e+05	USNS Raymond\nFPO AE 09386
5	80175.75416	4.988408	6.104512	4.04	26748.42842	1.068138e+06	06039 Jennifer Islands Apt. 443\nTracyport, KS...
6	64698.46343	6.025336	8.147760	3.41	60828.24909	1.502056e+06	4759 Daniel Shoals Suite 442\nNguyenburgh, CO ...
7	78394.33928	6.989780	6.620478	2.42	36516.35897	1.573937e+06	972 Joyce Viaduct\nLake William, TN 17778-6483

info

In [4]:

```
# to identify missing values
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   Avg. Area Income    5000 non-null   float64 
 1   Avg. Area House Age 5000 non-null   float64 
 2   Avg. Area Number of Rooms 5000 non-null   float64 
 3   Avg. Area Number of Bedrooms 5000 non-null   float64 
 4   Area Population     5000 non-null   float64 
 5   Price               5000 non-null   float64 
 6   Address             5000 non-null   object  
dtypes: float64(6), object(1)
memory usage: 273.6+ KB
```

describe

In [5]:

```
# to display summary of the dataset
data.describe()
```

Out[5]:

	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%	61480.562390	5.322283	6.299250	3.140000	29403.928700	9.975771e+05
50%	68804.286405	5.970429	7.002902	4.050000	36199.406690	1.232669e+06
75%	75783.338665	6.650808	7.665871	4.490000	42861.290770	1.471210e+06
max	107701.748400	9.519088	10.759588	6.500000	69621.713380	2.469066e+06

columns

In [6]:

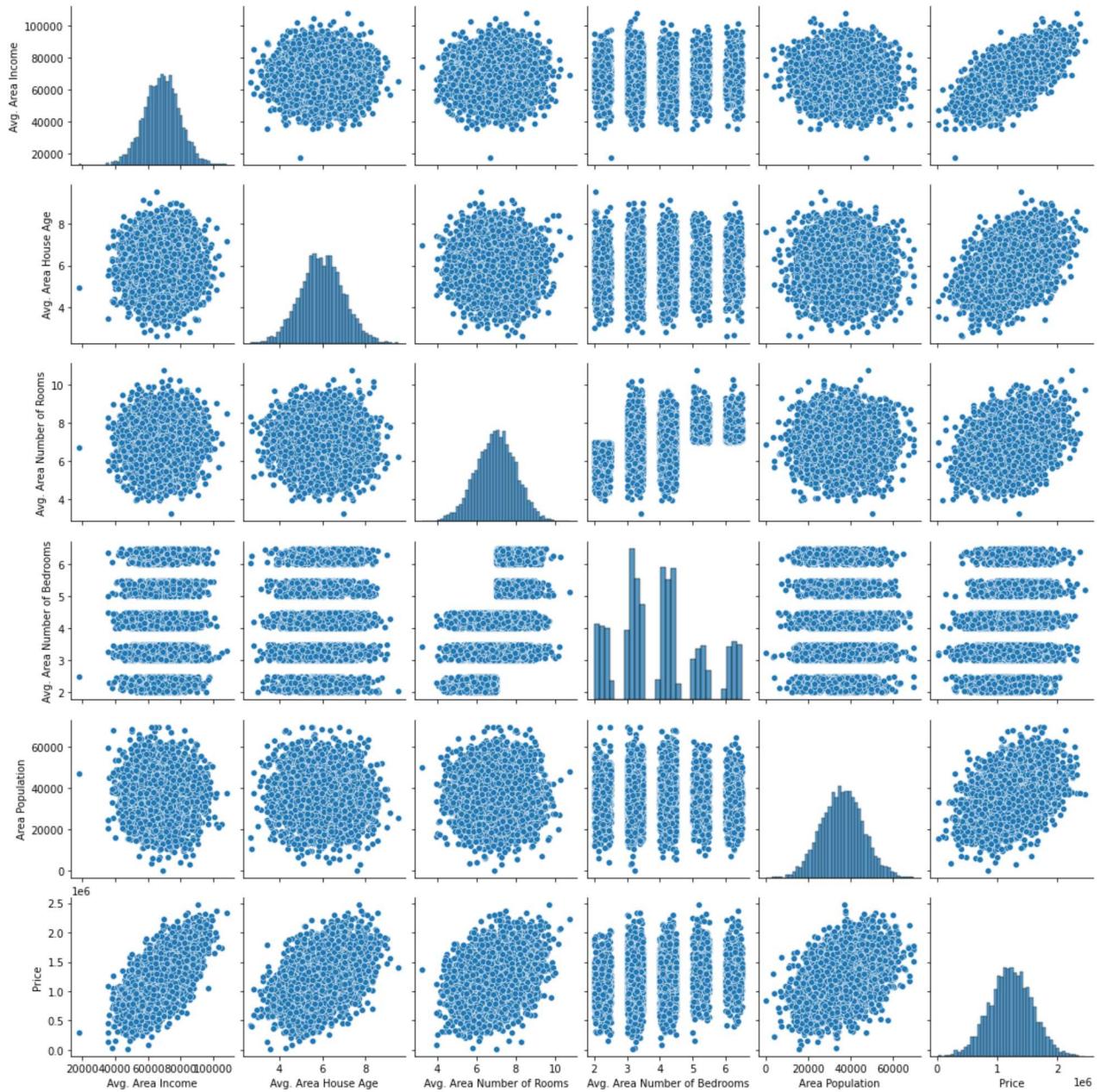
```
# to display headings of the dataset
data.columns
```

```
Out[6]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
       'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],
       dtype='object')
```

EDA and Visualization

```
In [7]: sns.pairplot(data)
```

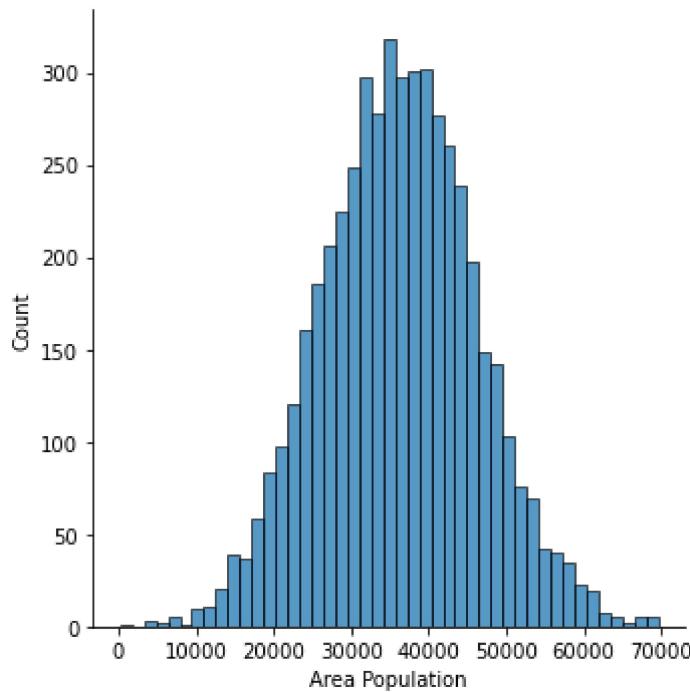
```
Out[7]: <seaborn.axisgrid.PairGrid at 0x1d0e1a8fb50>
```



distribution plot

```
In [8]: sns.displot(data["Area Population"])
```

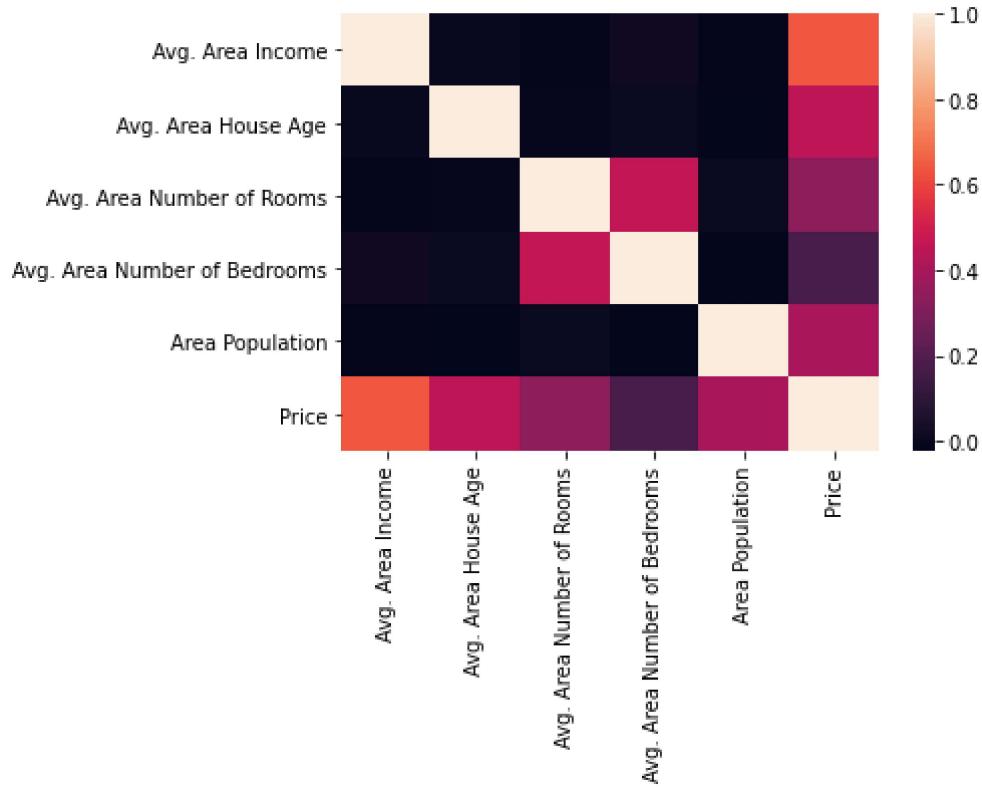
```
Out[8]: <seaborn.axisgrid.FacetGrid at 0x1d0e40407c0>
```



correlation

```
In [9]: da=data[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
           'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address']]
sns.heatmap(da.corr())
```

```
Out[9]: <AxesSubplot:>
```



To train the model-Model Building

we are going to train Linear Regression model; we need to split out data into two variables x and y where x is independent variable and y is dependent variable on x (output) we could ignore address column as it is not required for our model

```
In [10]: x=da[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
      'Avg. Area Number of Bedrooms', 'Area Population']]
y=da['Price']
```

```
In [11]: # to split my dataset into training and test data
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [12]: from sklearn.linear_model import LinearRegression
lr= LinearRegression()
lr.fit(x_train,y_train)
```

```
Out[12]: LinearRegression()
```

```
In [13]: print(lr.intercept_)
```

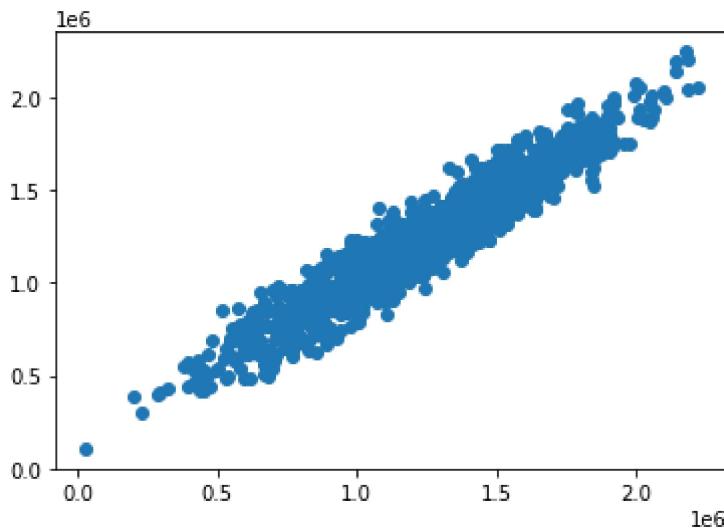
```
-2631208.336222648
```

```
In [14]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

	Co-efficient
Avg. Area Income	21.586568
Avg. Area House Age	164619.146776
Avg. Area Number of Rooms	121148.527381
Avg. Area Number of Bedrooms	890.502703
Area Population	15.162119

```
In [15]: prediction=lr.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[15]: <matplotlib.collections.PathCollection at 0x1d0e639ac10>
```



```
In [16]: print(lr.score(x_test,y_test))
```

```
0.9097608954127643
```

```
In [17]: lr.score(x_train,y_train)
```

```
Out[17]: 0.9211322953007236
```

Ridge regression

```
In [18]: from sklearn.linear_model import Ridge,Lasso
```

```
In [19]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
rr.score(x_test,y_test)
```

```
Out[19]: 0.9097393233157853
```

```
In [20]: rr.score(x_train,y_train)
```

```
Out[20]: 0.9211293671662266
```

LaSSO regression

```
In [21]: la=Lasso(alpha=10)
la.fit(x_train,y_train)
la.score(x_train,y_train)
```

```
Out[21]: 0.921132293635673
```

```
In [22]: la.score(x_test,y_test)
```

```
Out[22]: 0.9097601510696592
```

```
In [23]: from sklearn.linear_model import ElasticNet  
en=ElasticNet()  
en.fit(x_train,y_train)
```

```
Out[23]: ElasticNet()
```

```
In [24]: print(en.coef_)
```

```
[2.16085718e+01 1.08511339e+05 7.61870113e+04 1.40360683e+04  
1.52026864e+01]
```

```
In [25]: print(en.intercept_)
```

```
-2037832.090684975
```

```
In [26]: predict=en.predict(x_test)
```

```
In [27]: print(en.score(x_test,y_test))
```

```
0.8665434511757972
```

```
In [28]: from sklearn import metrics
```

```
In [29]: print("Mean Absolute error:",metrics.mean_absolute_error(y_test,predict))
```

```
Mean Absolute error: 99982.91001381187
```

```
In [30]: print("Mean Squared error:",metrics.mean_squared_error(y_test,predict))
```

```
Mean Squared error: 15289753787.911066
```

```
In [31]: print("Root squared error:",np.sqrt(metrics.mean_squared_error(y_test,predict)))
```

```
Root squared error: 123651.74397440202
```

```
In [ ]:
```