

# 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 him to estimate of what the house sell for.

```
In [22]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
In [3]: df = pd.read_csv("10_USA_Housing.csv")
df
```

Out[3]:

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanieltown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386
...	...	...	...	...	...	...	...
4995	60567.944140	7.830362	6.137356	3.46	22837.361035	1.060194e+06	USNS Williams\nFPO AP 30153-7653
4996	78491.275435	6.999135	6.576763	4.02	25616.115489	1.482618e+06	PSC 9258, Box 8489\nAPO AA 42991- 3352
4997	63390.686886	7.250591	4.805081	2.13	33266.145490	1.030730e+06	4215 Tracy Garden Suite 076\nJoshualand, VA 01...
4998	68001.331235	5.534388	7.130144	5.44	42625.620156	1.198657e+06	USS Wallace\nFPO AE 73316
4999	65510.581804	5.992305	6.792336	4.07	46501.283803	1.298950e+06	37778 George Ridges Apt. 509\nEast Holly,

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
							NV 2...

5000 rows × 7 columns

In [4]:

df.head()

Out[4]:

	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. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386

## Data cleaning and pre processing

In [5]:

df.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
```

In [6]:

df.describe()

Out[6]:

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

In [8]:

```
df.columns
```

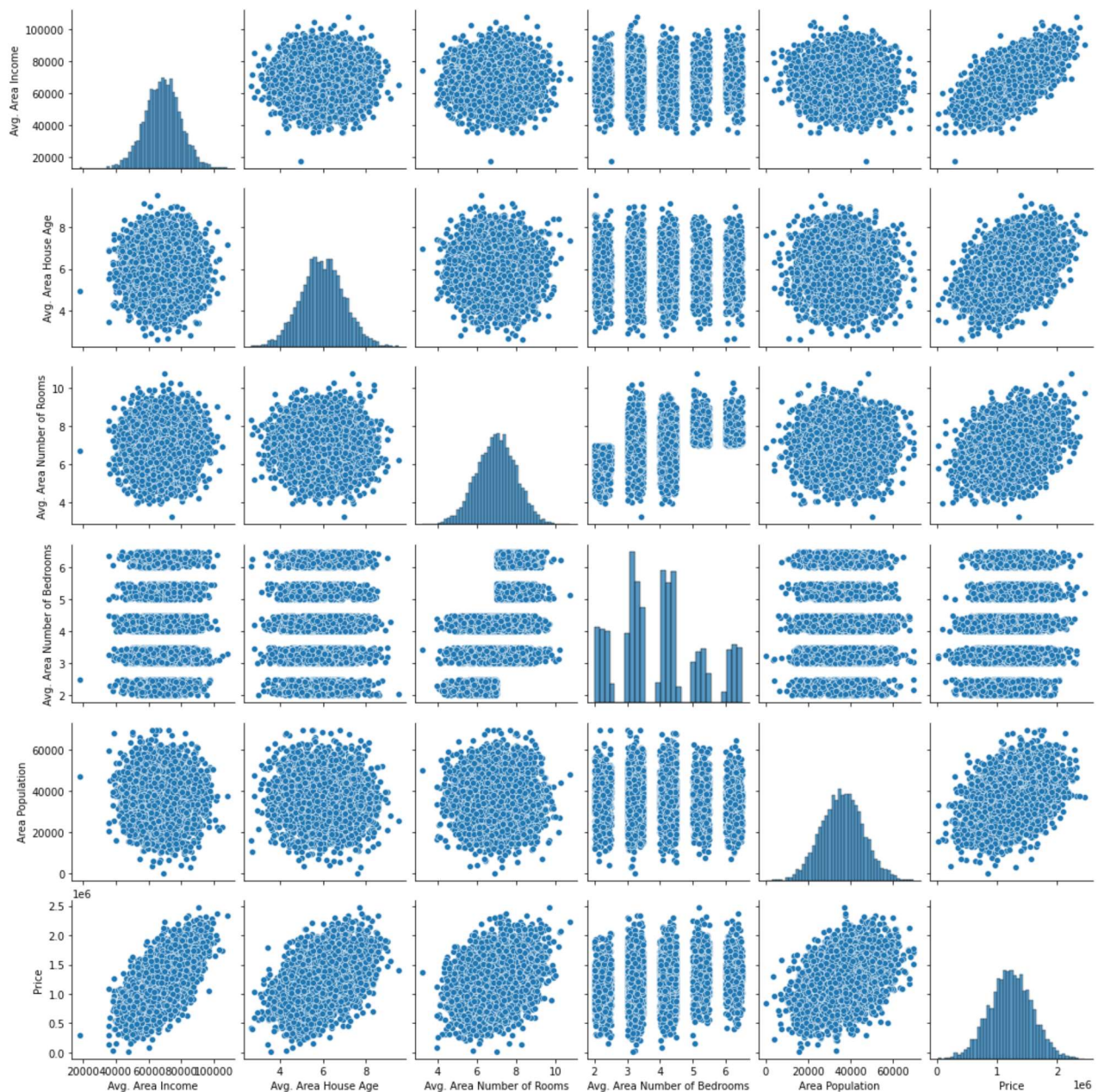
```
Out[8]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
              'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],
              dtype='object')
```

## EDA and VISUALIZATION

In [9]:

```
sns.pairplot(df)
```

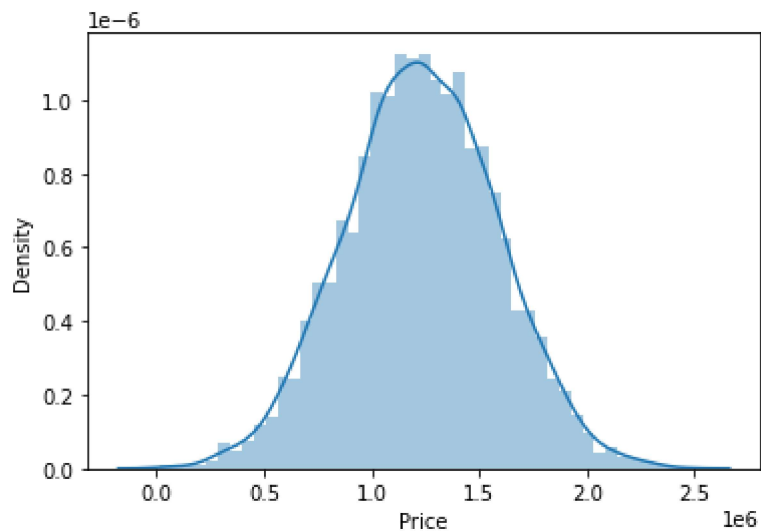
```
Out[9]: <seaborn.axisgrid.PairGrid at 0x190db374820>
```



```
In [10]: sns.distplot(df["Price"])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

```
Out[10]: <AxesSubplot:xlabel='Price', ylabel='Density'>
```



```
In [12]: df1 = df[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
                'Avg. Area Number of Bedrooms', 'Area Population', 'Price']]
```

```
In [14]: sns.heatmap(df1.corr())
```

Out[14]: <AxesSubplot:>



## To train the model - model building

going to train linear Regression model

- first split the data into two variables x and y where x is independent variable (input) and y is dependent on x (output)
- we can ignore address column as it is not required for our model

```
In [15]: x = df1[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
            'Avg. Area Number of Bedrooms', 'Area Population']]
        y = df1['Price']
```

## split the data into training and test data

```
In [19]: x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.3)
```

```
In [24]: lr = LinearRegression()
        lr.fit(x_train, y_train)
```

```
Out[24]: LinearRegression()
```

```
In [26]: lr.intercept_
```

```
Out[26]: -2639713.984629445
```

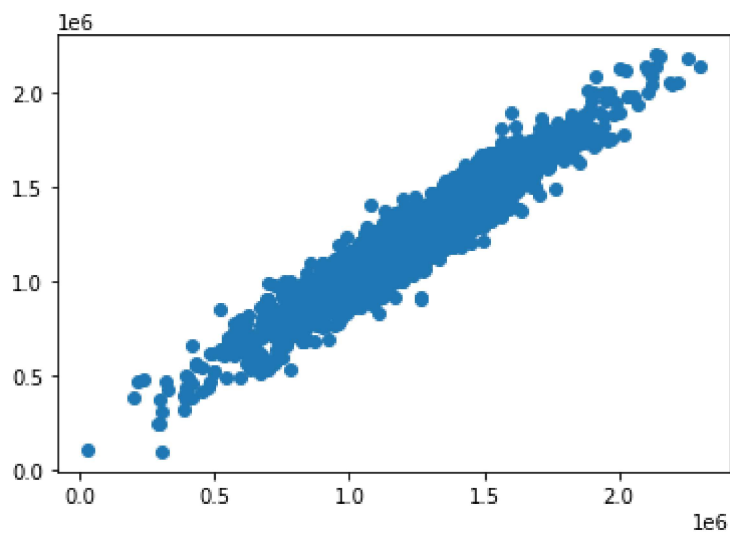
```
In [27]: coeff = pd.DataFrame(lr.coef_, x.columns, columns = ['Co-efficient'])
        coeff
```

```
Out[27]:
```

	Co-efficient
<b>Avg. Area Income</b>	21.577210
<b>Avg. Area House Age</b>	165324.413148
<b>Avg. Area Number of Rooms</b>	121051.833653
<b>Avg. Area Number of Bedrooms</b>	1227.694779
<b>Area Population</b>	15.308496

```
In [28]: prediction = lr.predict(x_test)
        plt.scatter(y_test, prediction)
```

```
Out[28]: <matplotlib.collections.PathCollection at 0x190ded61a90>
```



```
In [29]: lr.score(x_test,y_test)
```

```
Out[29]: 0.9164989189516786
```