# Problem Statement

A real estate agent want to 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 it.

# Linear Regression

# Import Libraries

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

In [3]:
```python
a=pd.read_csv("house.csv")
a
```

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 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 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 |
| 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, |

| 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

# To display top 10 rows

In [4]:
```python
a.head(10)
```

Out[4]:

| | Avg. Area Income | Avg. Area House Age | Avg. Area Number of Rooms | Avg. Area Number of Bedrooms | Area Population | Price | Address |
|---|---|---|---|---|---|---|---|
| 0 | 79545.45857 | 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 |
| 8 | 59927.66081 | 5.362126 | 6.393121 | 2.30 | 29387.39600 | 7.988695e+05 | USS Gilbert\nFPO AA 20957 |
| 9 | 81885.92718 | 4.423672 | 8.167688 | 6.10 | 40149.96575 | 1.545155e+06 | Unit 9446 Box 0958\nDPO AE 97025 |

# To find Missing values

In [5]:
```python
a.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
```

# To display summary of statistics

In [7]:
```python
a.describe()
```

Out[7]:

|  | 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 |

# To display column heading

In [8]:
```python
a.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'],
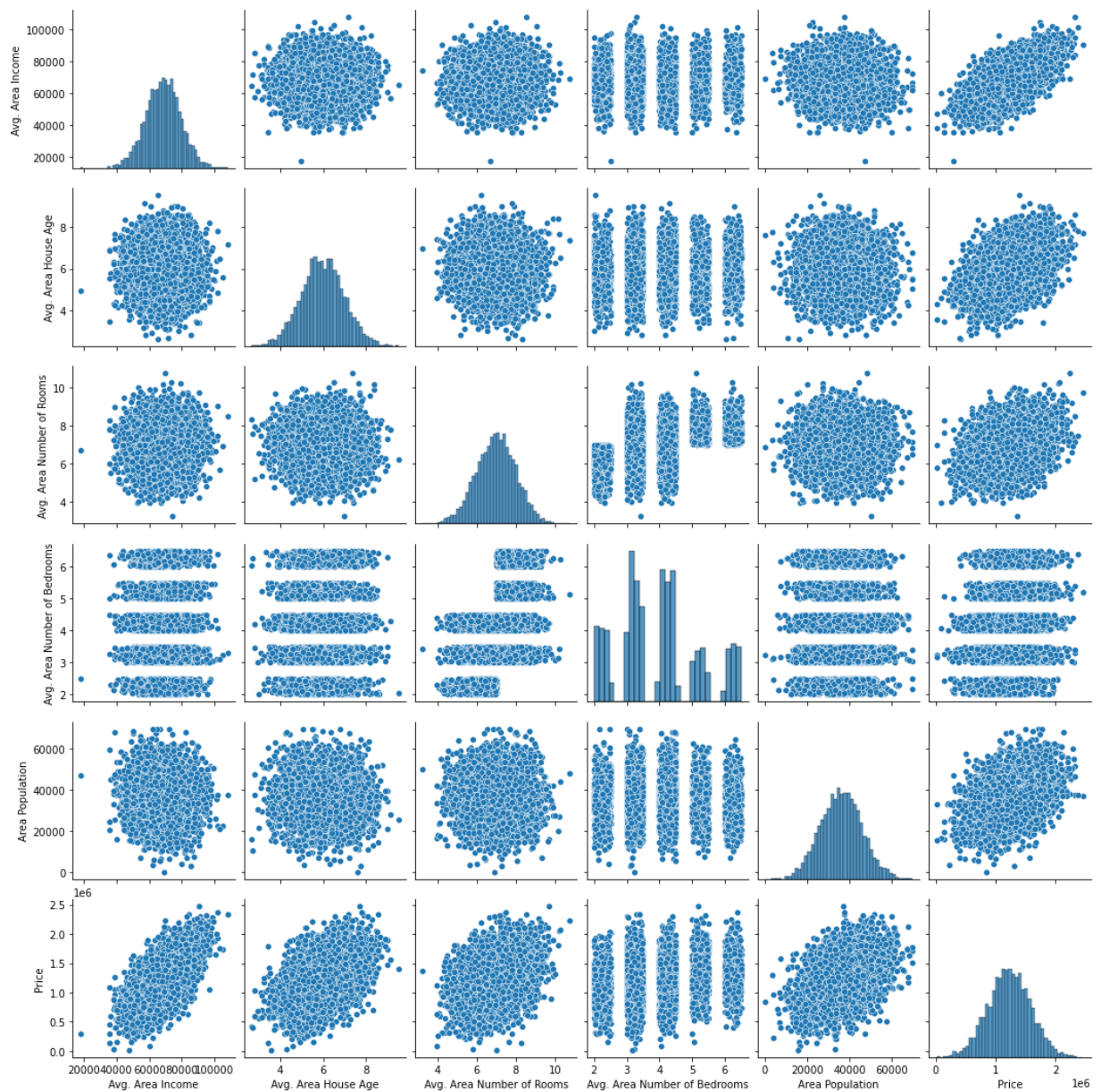      dtype='object')

# Pairplot

In [9]:
```python
sns.pairplot(a)
```

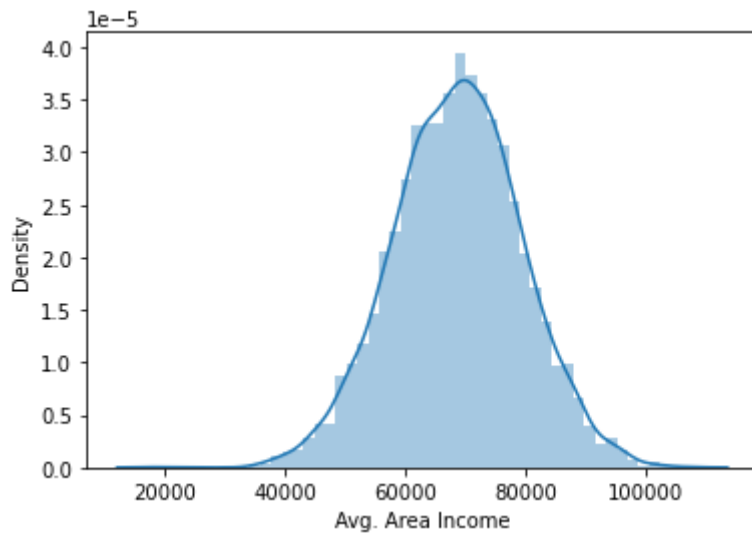Out[9]: <seaborn.axisgrid.PairGrid at 0x20472ec8ac0>

# Distribution Plot

```
In [12]:   sns.distplot(a['Avg. Area Income'])
```

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[12]:   <AxesSubplot:xlabel='Avg. Area Income', ylabel='Density'>

# Correlation

```
In [13]:   b=a[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
                 'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address']]
           sns.heatmap(b.corr())
```

Out[13]:   <AxesSubplot:>



# Train the model - Model Building

We are going to train linear regression model: We need to split out data into 2 variables x,y
where x is independant and y is dependant on x(output). We could ignore address column as it
is not required for our model.

In [25]:
```python
g=b[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
     'Avg. Area Number of Bedrooms', 'Area Population']]
h=b['Price']
```

# To split dataset into training end test

In [26]:
```python
from sklearn.model_selection import train_test_split
g_train,g_test,h_train,h_test=train_test_split(g,h,test_size=0.5)
```

# To run the model

In [20]:
```python
from sklearn.linear_model import LinearRegression
```

In [27]:
```python
lr=LinearRegression()
lr.fit(g_train,h_train)
```

Out[27]: LinearRegression()

In [29]:
```python
print(lr.intercept_)
```

-2656178.1464716895

# Coeffecient

In [32]:
```python
coeff=pd.DataFrame(lr.coef_,g.columns,columns=['Co-effecient'])
coeff
```
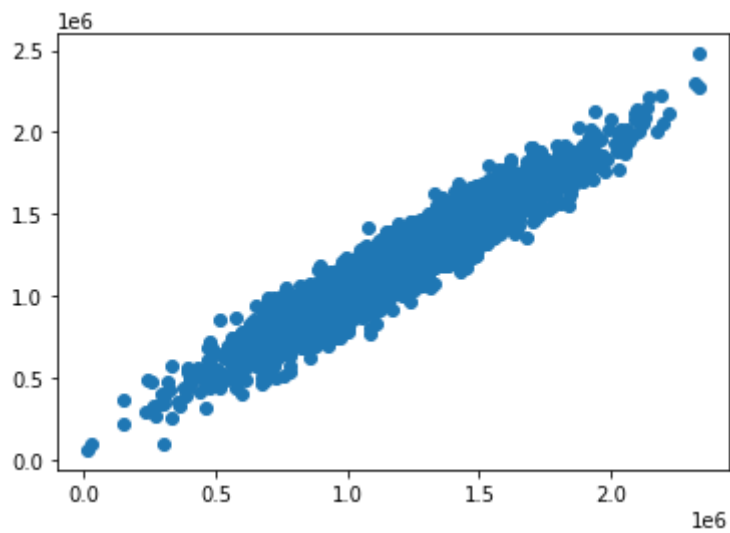
Out[32]:

|  | Co-effecient |
| --- | --- |
| Avg. Area Income | 21.623153 |
| Avg. Area House Age | 167284.467720 |
| Avg. Area Number of Rooms | 121650.958678 |
| Avg. Area Number of Bedrooms | 1261.418030 |
| Area Population | 15.203649 |

# Best Fit line

In [34]:
```python
prediction=lr.predict(g_test)
plt.scatter(h_test,prediction)
```

Out[34]: <matplotlib.collections.PathCollection at 0x20477b97b80>

# To find score

In [35]:
```python
print(lr.score(g_test,h_test))
```

0.9122114121025614

In [ ]: