

Lecture 9- Part 2

Linear Regression

Import data

- The data contains the following columns:
- 'Avg. Area Income': Avg. Income of residents of the city house is located in.
- 'Avg. Area House Age': Avg Age of Houses in same city
- 'Avg. Area Number of Rooms': Avg Number of Rooms for Houses in same city
- 'Avg. Area Number of Bedrooms': Avg Number of Bedrooms for Houses in same city
- 'Area Population': Population of city house is located in
- 'Price': Price that the house sold at
- 'Address': Address for the house

```
In [2]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 %matplotlib inline
```

```
In [30]: 1 data = pd.read_csv('Downloads/House-price.csv')
        2 data.head(5)
```

Out[30]:

	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 A 674\nLaurabury, N 3701
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson View Suite 079\nLal Kathleen, CA
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabe Stravenue\nDanieltow WI 06482
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO # 448;
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFF AE 093

In [31]: 1 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
```

In [32]: 1 data.describe()

Out[32]:

	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 [33]: 1 data.columns

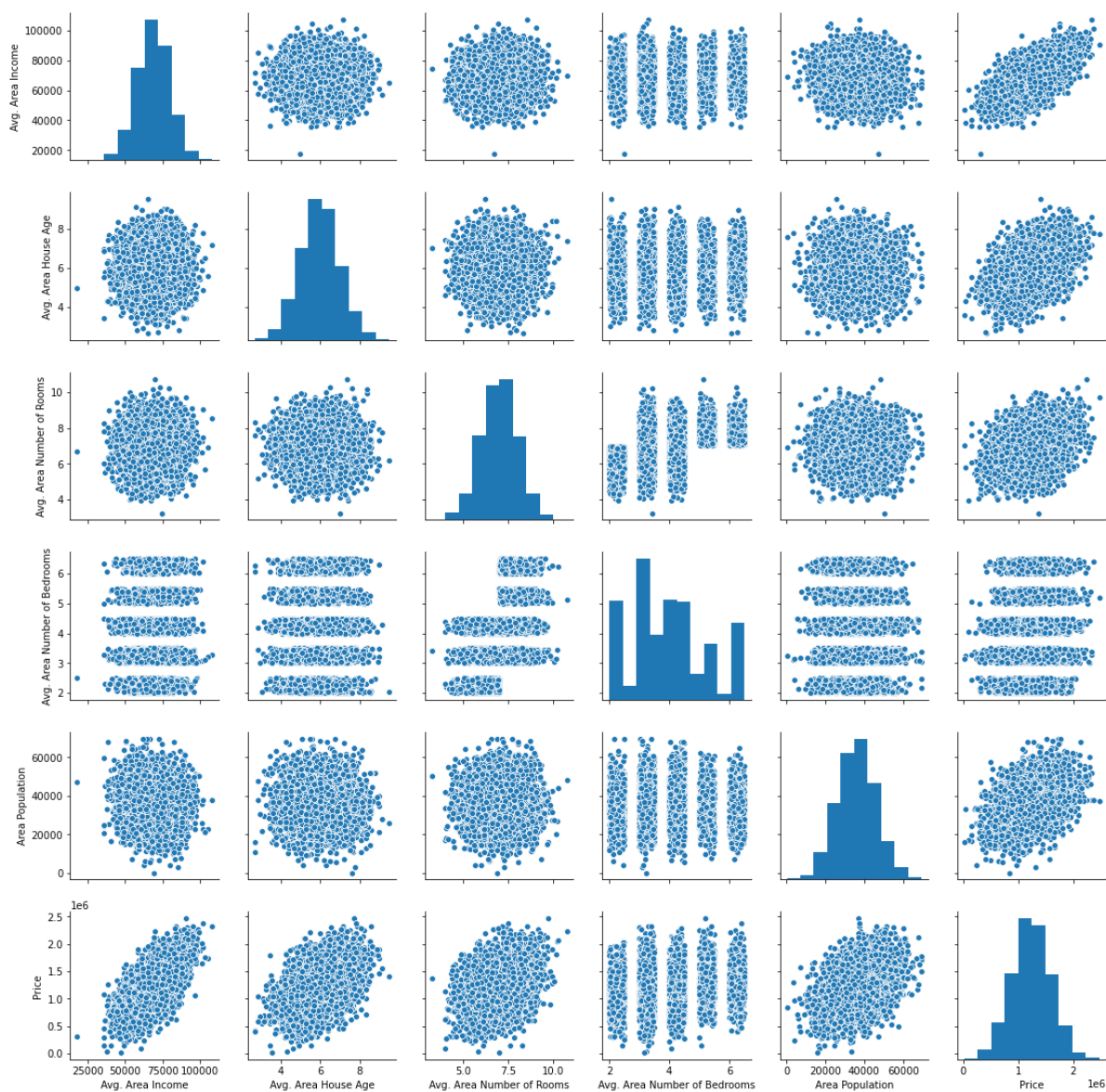
Out[33]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms', 'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'], dtype='object')

Analysis

Let's create some simple plots to check out the data!

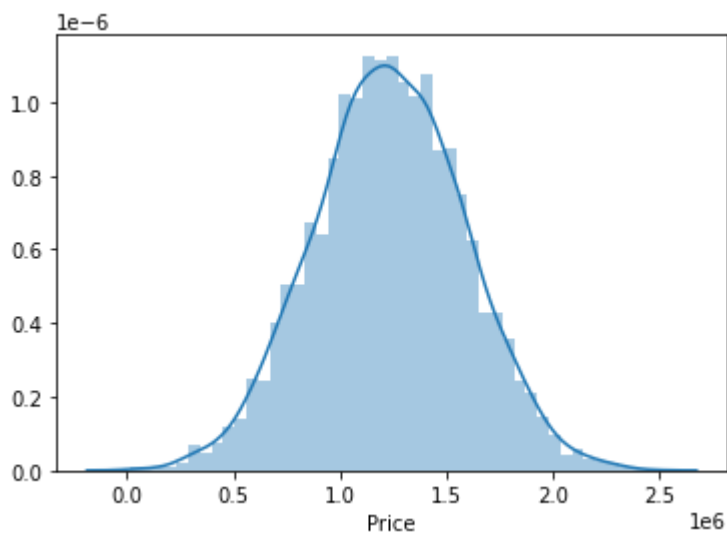
```
In [34]: 1 sns.pairplot(data)
```

```
Out[34]: <seaborn.axisgrid.PairGrid at 0x11632bcd708>
```



```
In [37]: 1 sns.distplot(data['Price'])
```

```
Out[37]: <matplotlib.axes._subplots.AxesSubplot at 0x11633093f48>
```



```
In [38]: 1 sns.heatmap(data.corr())
```

```
Out[38]: <matplotlib.axes._subplots.AxesSubplot at 0x11633d1dec8>
```



Training a Linear Regression Model

Let's now begin to train our regression model! We will need to first split up our data into an X array that contains the features to train on, and a y array with the target variable, in this case the Price column. We will toss out the Address column because it only has text info that the linear

regression model can't use.

X and y arrays

```
In [41]: 1 X = data[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of  
2           'Avg. Area Number of Bedrooms', 'Area Population']]  
3 y = data['Price']
```

Train Test Split

Now let's split the data into a training set and a testing set. We will train our model on the training set and then use the test set to evaluate the model.

```
In [42]: 1 from sklearn.model_selection import train_test_split
```

```
In [43]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, r
```

Creating and Training the Model

```
In [44]: 1 from sklearn.linear_model import LinearRegression
```

```
In [45]: 1 lm = LinearRegression()
```

```
In [46]: 1 lm.fit(X_train,y_train)
```

```
Out[46]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

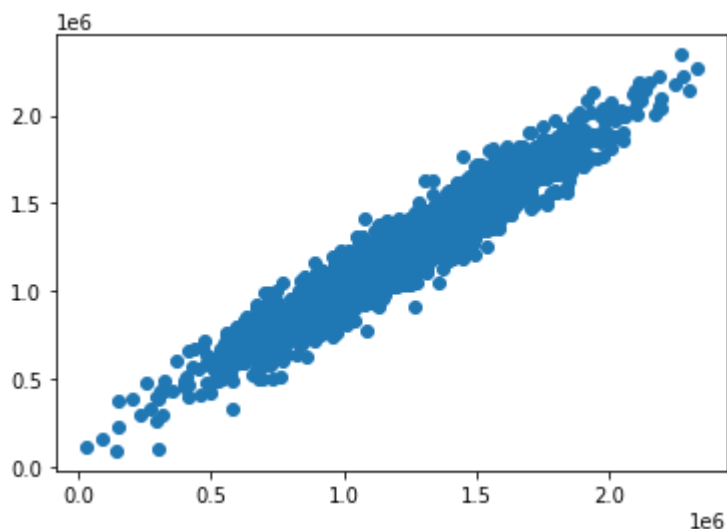
Predictions from our Model

Let's grab predictions off our test set and see how well it did!

```
In [47]: 1 predictions = lm.predict(X_test)
```

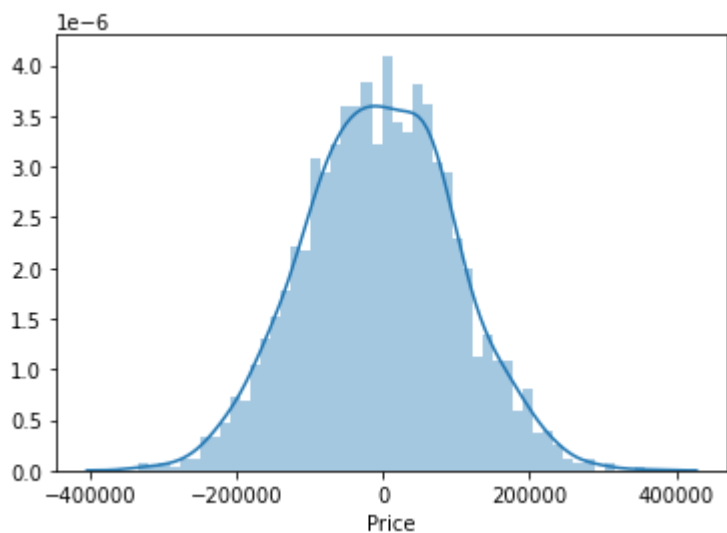
```
In [48]: 1 plt.scatter(y_test,predictions)
```

```
Out[48]: <matplotlib.collections.PathCollection at 0x1163542e9c8>
```



Residual Histogram

```
In [49]: 1 sns.distplot((y_test-predictions),bins=50);
```



Regression Evaluation Metrics

Here are three common evaluation metrics for regression problems:

Mean Absolute Error (MAE) is the mean of the absolute value of the errors:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE) is the mean of the squared errors:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Comparing these metrics:

- **MAE** is the easiest to understand, because it's the average error.
- **MSE** is more popular than MAE, because MSE "punishes" larger errors, which tends to be useful in the real world.
- **RMSE** is even more popular than MSE, because RMSE is interpretable in the "y" units.

All of these are loss functions, because we want to minimize them

```
In [50]: 1 from sklearn import metrics
```

```
In [51]: 1 print('MAE:', metrics.mean_absolute_error(y_test, predictions))
2 print('MSE:', metrics.mean_squared_error(y_test, predictions))
3 print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, predictions)))
```

```
MAE: 82288.22251914957
MSE: 10460958907.209507
RMSE: 102278.82922291156
```

```
In [ ]: 1
```