

# Linear Regression (US Housing)

October 26, 2016

\_ # Linear Regression with Python

## 0.1 U.S Housing Project

- '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]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
In [3]: USAhousing = pd.read_csv('USA_Housing.csv')
```

```
In [4]: USAhousing.head()
```

```
Out[4]:
```

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	
0	79545.458574	5.682861	7.009188	
1	79248.642455	6.002900	6.730821	
2	61287.067179	5.865890	8.512727	
3	63345.240046	7.188236	5.586729	
4	59982.197226	5.040555	7.839388	

	Avg. Area Number of Bedrooms	Area Population	Price	
0	4.09	23086.800503	1.059034e+06	
1	3.09	40173.072174	1.505891e+06	
2	5.13	36882.159400	1.058988e+06	
3	3.26	34310.242831	1.260617e+06	
4	4.23	26354.109472	6.309435e+05	

	Address
0	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	USS Barnett\nFPO AP 44820
4	USNS Raymond\nFPO AE 09386

```
In [5]: USAhousing.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
Avg. Area Income                5000 non-null float64
Avg. Area House Age             5000 non-null float64
Avg. Area Number of Rooms       5000 non-null float64
Avg. Area Number of Bedrooms    5000 non-null float64
Area Population                 5000 non-null float64
Price                          5000 non-null float64
Address                        5000 non-null object
dtypes: float64(6), object(1)
memory usage: 273.5+ KB

```

```
In [6]: USAhousing.describe()
```

```

Out[6]:
      Avg. Area Income  Avg. Area House Age  Avg. Area Number of Rooms \
count      5000.000000          5000.000000          5000.000000
mean      68583.108984           5.977222           6.987792
std       10657.991214           0.991456           1.005833
min       17796.631190           2.644304           3.236194
25%       61480.562388           5.322283           6.299250
50%       68804.286404           5.970429           7.002902
75%       75783.338666           6.650808           7.665871
max       107701.748378           9.519088          10.759588

      Avg. Area Number of Bedrooms  Area Population      Price
count          5000.000000          5000.000000  5.000000e+03
mean              3.981330          36163.516039  1.232073e+06
std              1.234137           9925.650114  3.531176e+05
min              2.000000           172.610686  1.593866e+04
25%              3.140000          29403.928702  9.975771e+05
50%              4.050000          36199.406689  1.232669e+06
75%              4.490000          42861.290769  1.471210e+06
max              6.500000          69621.713378  2.469066e+06

```

```
In [7]: USAhousing.columns
```

```

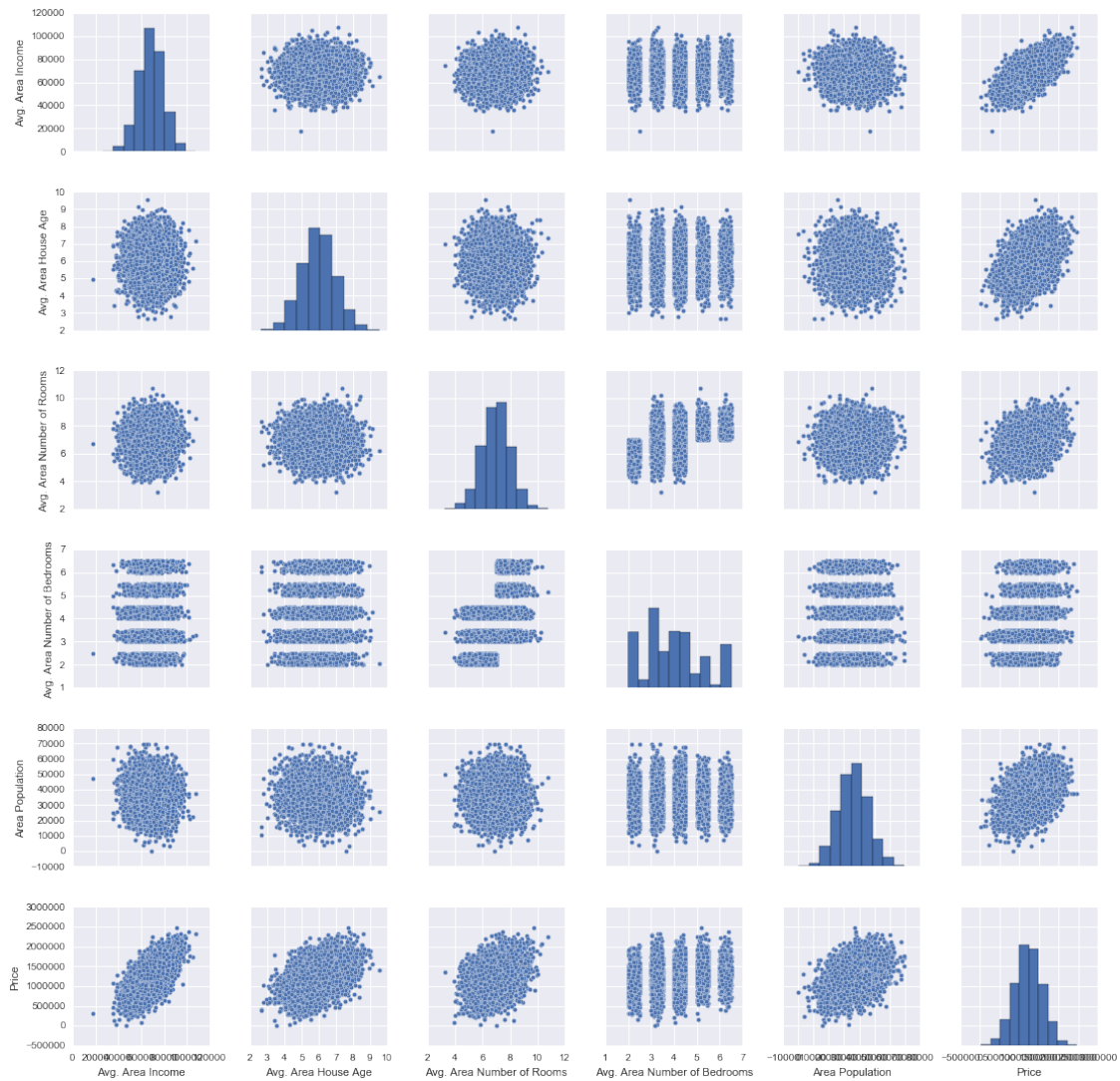
Out[7]: Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',
              'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],
              dtype='object')

```

## 1 EDA

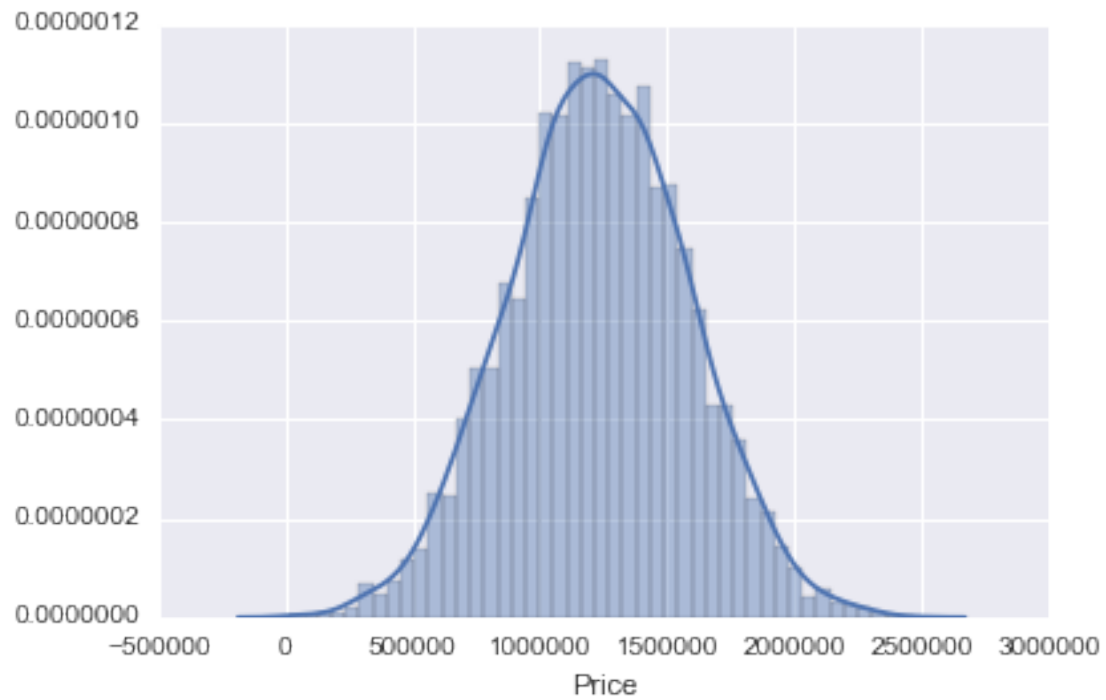
```
In [9]: sns.pairplot(USAhousing)
```

```
Out[9]: <seaborn.axisgrid.PairGrid at 0xd15dbe3da0>
```



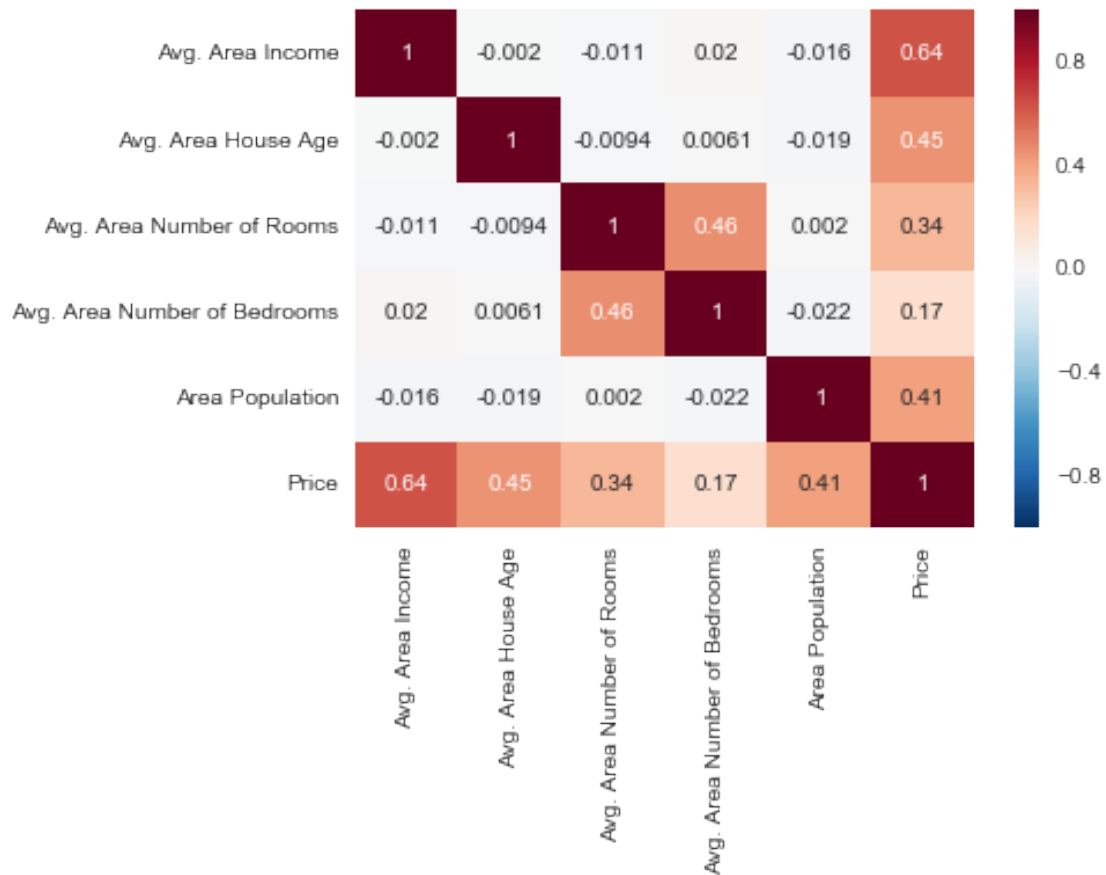
```
In [11]: sns.distplot(USAhousing['Price'])
```

```
Out[11]: <matplotlib.axes._subplots.AxesSubplot at 0xd1611c1278>
```



```
In [12]: sns.heatmap(USAhousing.corr(),annot=True)
```

```
Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0xd1612486a0>
```



## 1.1 Training a Linear Regression Model

In [ ]: Set your dependent variables for x and your dependent variable for Y

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

## 1.2 Train Test Split

```
In [17]: from sklearn.cross_validation import train_test_split
```

```
In [18]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=101)
```

## 1.3 Creating and Training the Model

```
In [19]: from sklearn.linear_model import LinearRegression
```

```
In [20]: lm = LinearRegression()
```

```
In [23]: lm.fit(X_train,y_train)
```

```
Out[23]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)
```

## 1.4 Model Evaluation

```
In [24]: # print the intercept  
print(lm.intercept_)
```

-2640159.79685

```
In [277]: coeff_df = pd.DataFrame(lm.coef_,X.columns,columns=['Coefficient'])  
coeff_df
```

```
Out[277]:
```

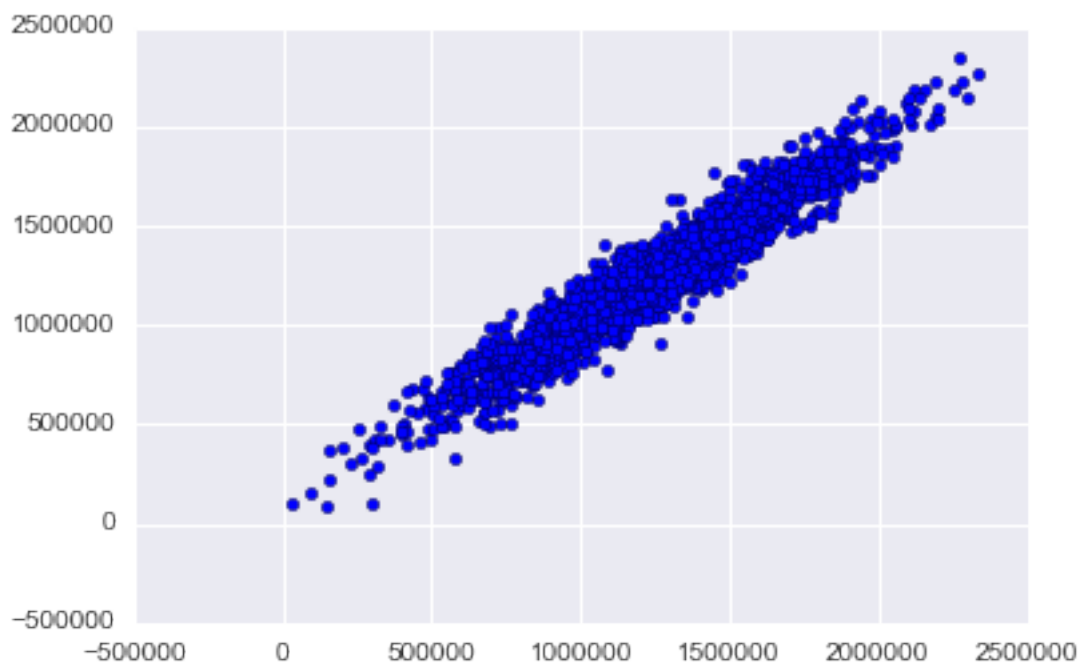
	Coefficient
Avg. Area Income	21.528276
Avg. Area House Age	164883.282027
Avg. Area Number of Rooms	122368.678027
Avg. Area Number of Bedrooms	2233.801864
Area Population	15.150420

## 1.5 Predictions from our Model

```
In [279]: predictions = lm.predict(X_test)
```

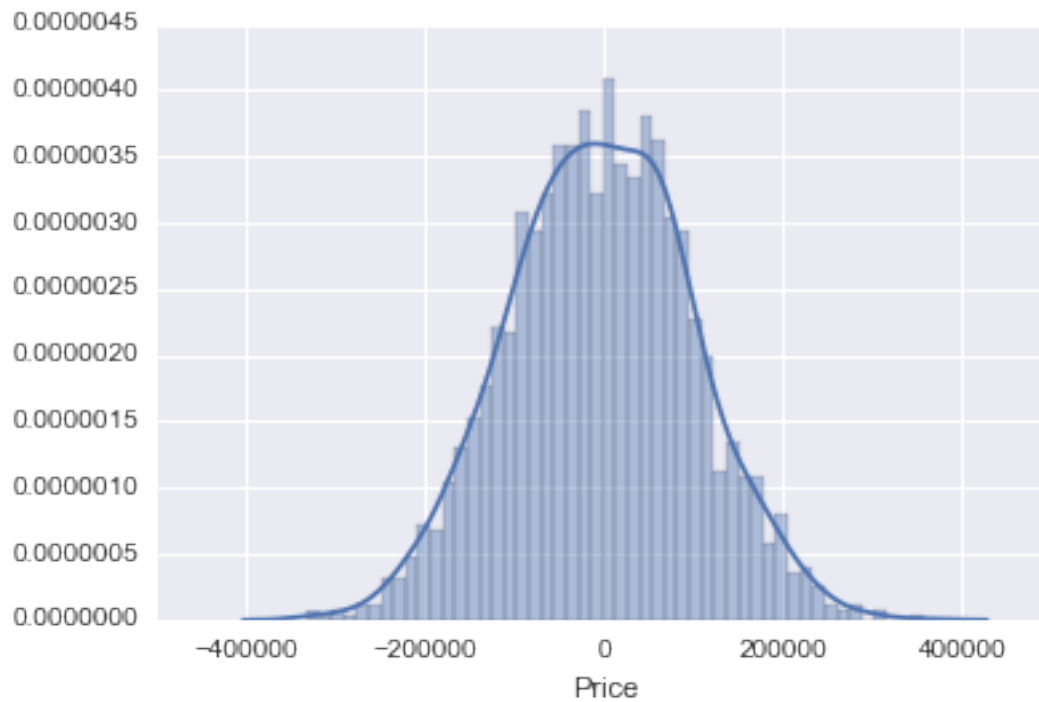
```
In [282]: plt.scatter(y_test,predictions)
```

```
Out[282]: <matplotlib.collections.PathCollection at 0x142622c88>
```



## Residual Histogram

```
In [281]: sns.distplot((y_test-predictions),bins=50);
```



```
In [275]: from sklearn import metrics
```

```
In [276]: print('MAE:', metrics.mean_absolute_error(y_test, predictions))  
          print('MSE:', metrics.mean_squared_error(y_test, predictions))  
          print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, predictions)))
```

```
MAE: 82288.2225191  
MSE: 10460958907.2  
RMSE: 102278.829223
```

## 1.6 $R^2$

```
In [26]: lm.score(X,y)
```

```
Out[26]: 0.91795587252010413
```

```
In [ ]:
```