

Model development

January 1, 2024

1 Hands-on practice lab: Model Development

2 Objectives

- Use Linear Regression in one variable to fit the parameters to a model
- Use Linear Regression in multiple variables to fit the parameters to a model
- Use Polynomial Regression in single variable to fit the parameters to a model
- Create a pipeline for performing linear regression using multiple features in polynomial scaling
- Evaluate the performance of different forms of regression on basis of MSE and R^2 parameters

3 Setup

For this lab, we will be using the following libraries:

- `pandas` for managing the data.
- `numpy` for mathematical operations.
- `sklearn` for machine learning and machine-learning-pipeline related functions.
- `seaborn` for visualizing the data.
- `matplotlib` for additional plotting tools.

3.0.1 Importing Required Libraries

We recommend you import all required libraries in one place (here):

```
[2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import StandardScaler, PolynomialFeatures
from sklearn.pipeline import Pipeline
from sklearn.metrics import mean_squared_error, r2_score
import warnings
warnings.filterwarnings("ignore", category=UserWarning)
%matplotlib inline
```

3.0.2 Importing the dataset

```
[4]: path = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/
↳IBMDeveloperSkillsNetwork-DA0101EN-Coursera/laptop_pricing_dataset_mod2.csv"
```

Load the dataset into a pandas dataframe

```
[6]: df = pd.read_csv("laptops.csv", header=0)
```

```
[7]: #https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/
↳IBMDeveloperSkillsNetwork-DA0101EN-Coursera/laptop_pricing_dataset_mod2.csv"
#df = pd.read_csv(filepath, header=None)
```

```
[8]: # show the first 5 rows using dataframe.head() method
print("The first 5 rows of the dataframe")
df.head(5)
```

The first 5 rows of the dataframe

```
[8]:
```

	Unnamed: 0.1	Unnamed: 0	Manufacturer	Category	GPU	OS	CPU_core	\
0	0	0	Acer	4	2	1	5	
1	1	1	Dell	3	1	1	3	
2	2	2	Dell	3	1	1	7	
3	3	3	Dell	4	2	1	5	
4	4	4	HP	4	2	1	7	

	Screen_Size_inch	CPU_frequency	RAM_GB	Storage_GB_SSD	Weight_pounds	\
0	14.0	0.551724	8	256	3.52800	
1	15.6	0.689655	4	256	4.85100	
2	15.6	0.931034	8	256	4.85100	
3	13.3	0.551724	8	128	2.69010	
4	15.6	0.620690	8	256	4.21155	

	Price	Price-binned	Screen-Full_HD	Screen-IPS_panel
0	978	Low	0	1
1	634	Low	1	0
2	946	Low	1	0
3	1244	Low	0	1
4	837	Low	1	0

4 Task 1 : Single Linear Regression

You have learnt that “CPU_frequency” is the parameter with the lowest p-value among the different features of the dataset. Create a single feature Linear Regression model that fits the pair of “CPU_frequency” and “Price” to find the model for prediction.

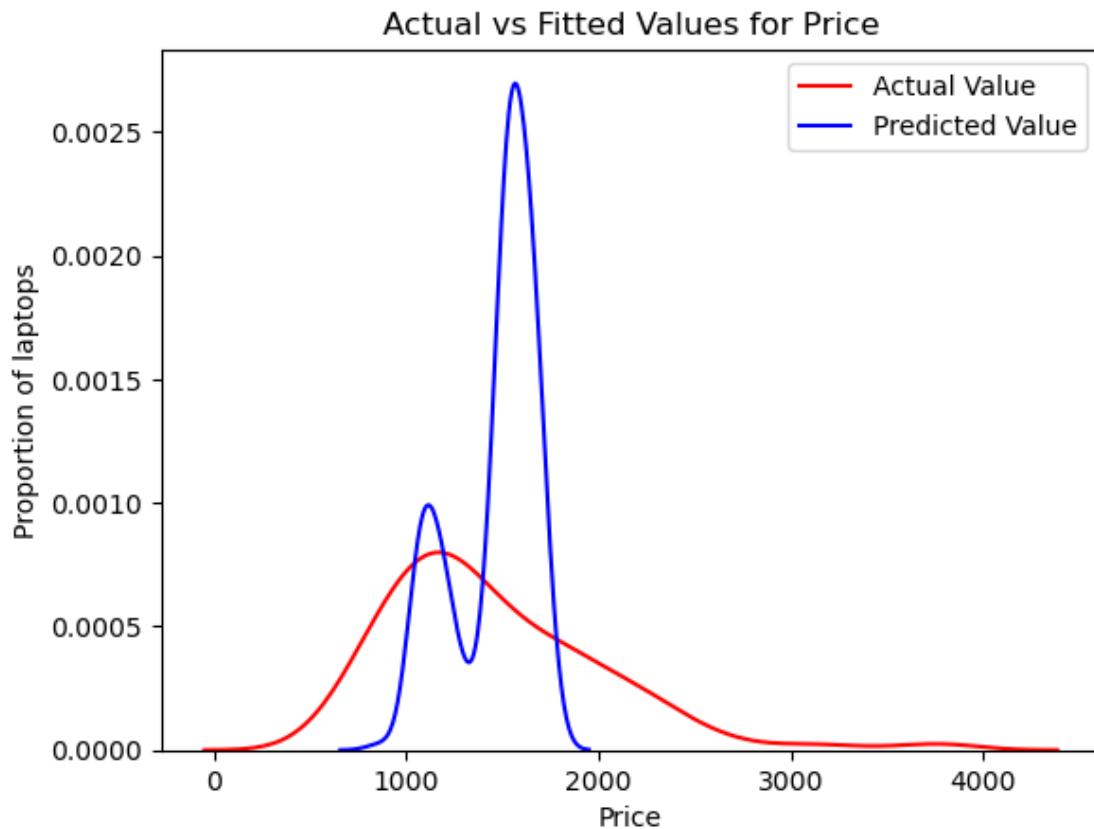
```
[13]: # Write your code below and press Shift+Enter to execute
lm=LinearRegression()
```

```
x=df[['CPU_frequency']]
y=df['Price']
lm.fit(x,y)
Yhat=lm.predict(x)
Yhat[0:5]
```

```
[13]: array([1073.07834392, 1277.93263722, 1636.42765051, 1073.07834392,
          1175.50549057])
```

Generate the Distribution plot for the predicted values and that of the actual values. How well did the model perform?

```
[15]: # Write your code below and press Shift+Enter to execute
ax1=sns.distplot(df['Price'],hist=False,color='r',label='Actual Values')
sns.distplot(Yhat,hist=False,color='b',label='Predicted Values')
plt.title('Actual vs Fitted Values for Price')
plt.xlabel('Price')
plt.ylabel('Proportion of laptops')
plt.legend(['Actual Value', 'Predicted Value'])
plt.show()
```



Evaluate the Mean Squared Error and R^2 score values for the model.

```
[ ]: # Write your code below and press Shift+Enter to execute
mse_slr = mean_squared_error(df['Price'], Yhat)
r2_score_slr = lm.score(X, Y)
print('The R-square for Linear Regression is: ', r2_score_slr)
print('The mean square error of price and predicted value is: ', mse_slr)
```

5 Task 2 - Multiple Linear Regression

The parameters which have a low enough p-value so as to indicate strong relationship with the 'Price' value are 'CPU_frequency', 'RAM_GB', 'Storage_GB_SSD', 'CPU_core', 'OS', 'GPU' and 'Category'. Use all these variables to create a Multiple Linear Regression system.

```
[16]: # Write your code below and press Shift+Enter to execute
x=df[['CPU_frequency', 'RAM_GB', 'Storage_GB_SSD', 'CPU_core', 'OS',
      ↵ 'GPU', 'Category']]
y=df['Price']
lm.fit(x,y)
r=lm.predict(x)
r[0:5]
```

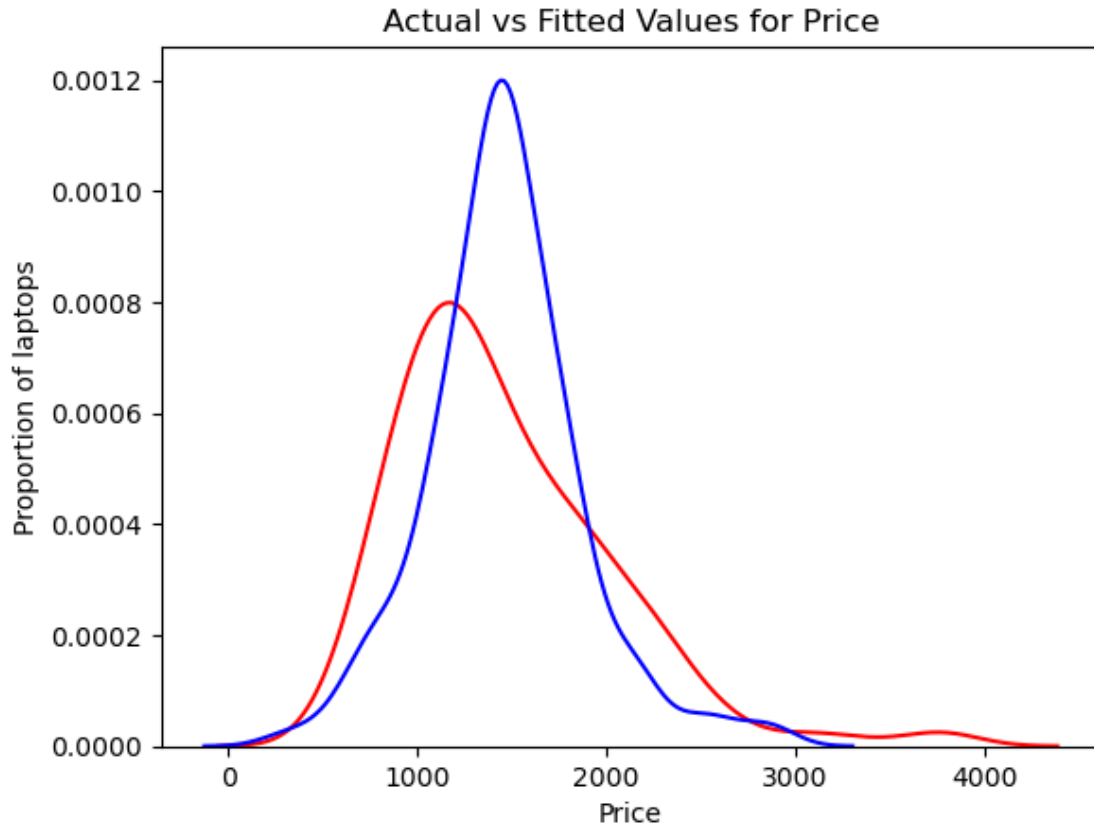
```
[16]: array([1345.51622771,  710.44905496, 1552.37242687, 1295.00681012,
          1543.13847022])
```

Plot the Distribution graph of the predicted values as well as the Actual values

```
[19]: # Write your code below and press Shift+Enter to execute
ax1 = sns.distplot(df['Price'], hist=False, color="r", label="Actual Value")
sns.distplot(r, hist=False, color="b", label="Fitted Values" , ax=ax1)

plt.title('Actual vs Fitted Values for Price')
plt.xlabel('Price')
plt.ylabel('Proportion of laptops')
```

```
[19]: Text(0, 0.5, 'Proportion of laptops')
```



Find the R^2 score and the MSE value for this fit. Is this better or worst than the performance of Single Linear Regression?

```
[21]: # Write your code below and press Shift+Enter to execute
mse_slr = mean_squared_error(df['Price'], r)
r2_score_slr = lm.score(x, y)
print('The R-square for Linear Regression is: ', r2_score_slr)
print('The mean square error of price and predicted value is: ', mse_slr)
```

The R-square for Linear Regression is: 0.5082509055187374

The mean square error of price and predicted value is: 161680.5726389311

6 Task 3 - Polynomial Regression

Use the variable “CPU_frequency” to create Polynomial features. Try this for 3 different values of polynomial degrees. Remember that polynomial fits are done using `numpy.polyfit`.

```
[24]: # Write your code below and press Shift+Enter to execute
X=df['CPU_frequency']
Y=df['Price']
X = X.to_numpy().flatten()
```

```
f1 = np.polyfit(X, Y, 1)
p1 = np.poly1d(f1)

f3 = np.polyfit(X, Y, 3)
p3 = np.poly1d(f3)

f5 = np.polyfit(X, Y, 5)
p5 = np.poly1d(f5)
```

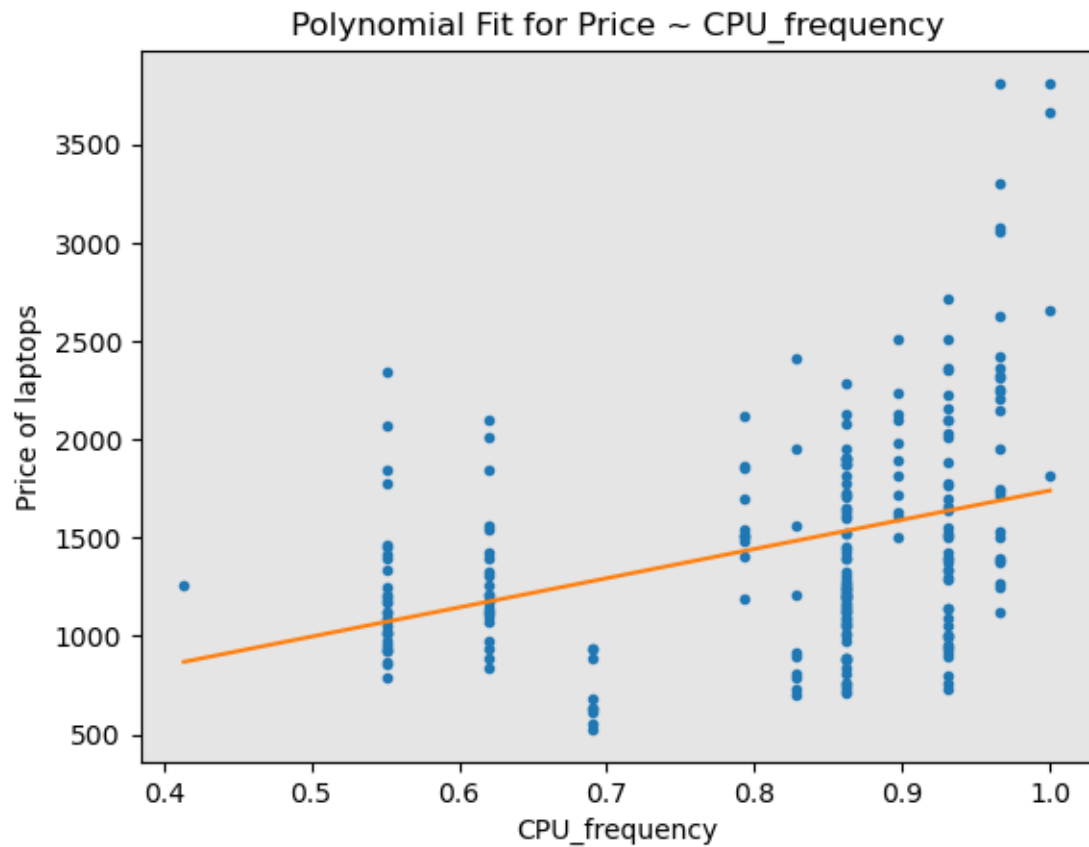
Plot the regression output against the actual data points to note how the data fits in each case. To plot the polynomial response over the actual data points, you have the function shown below.

```
[25]: def PlotPolly(model, independent_variable, dependent_variabble, Name):
        x_new = np.linspace(independent_variable.min(), independent_variable.
        ↪max(), 100)
        y_new = model(x_new)

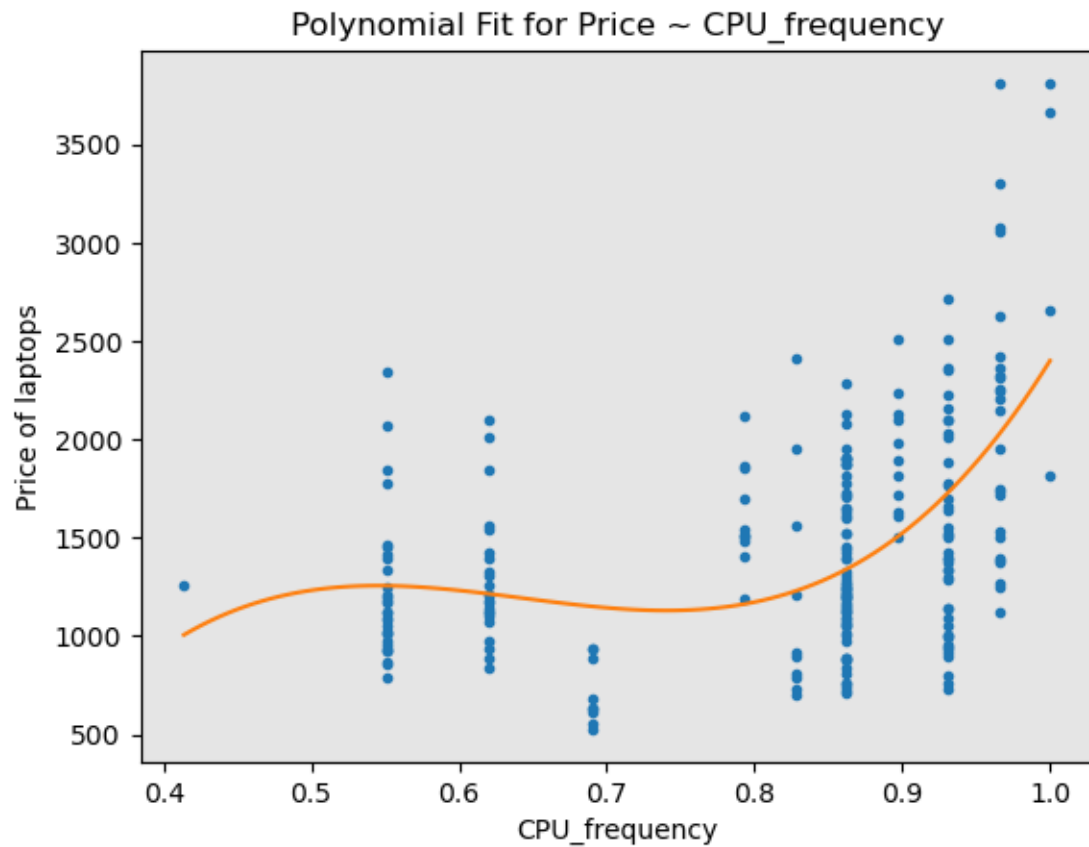
        plt.plot(independent_variable, dependent_variabble, '.', x_new, y_new, '-')
        plt.title(f'Polynomial Fit for Price ~ {Name}')
        ax = plt.gca()
        ax.set_facecolor((0.898, 0.898, 0.898))
        fig = plt.gcf()
        plt.xlabel(Name)
        plt.ylabel('Price of laptops')
```

Call this function for the 3 models created and get the required graphs.

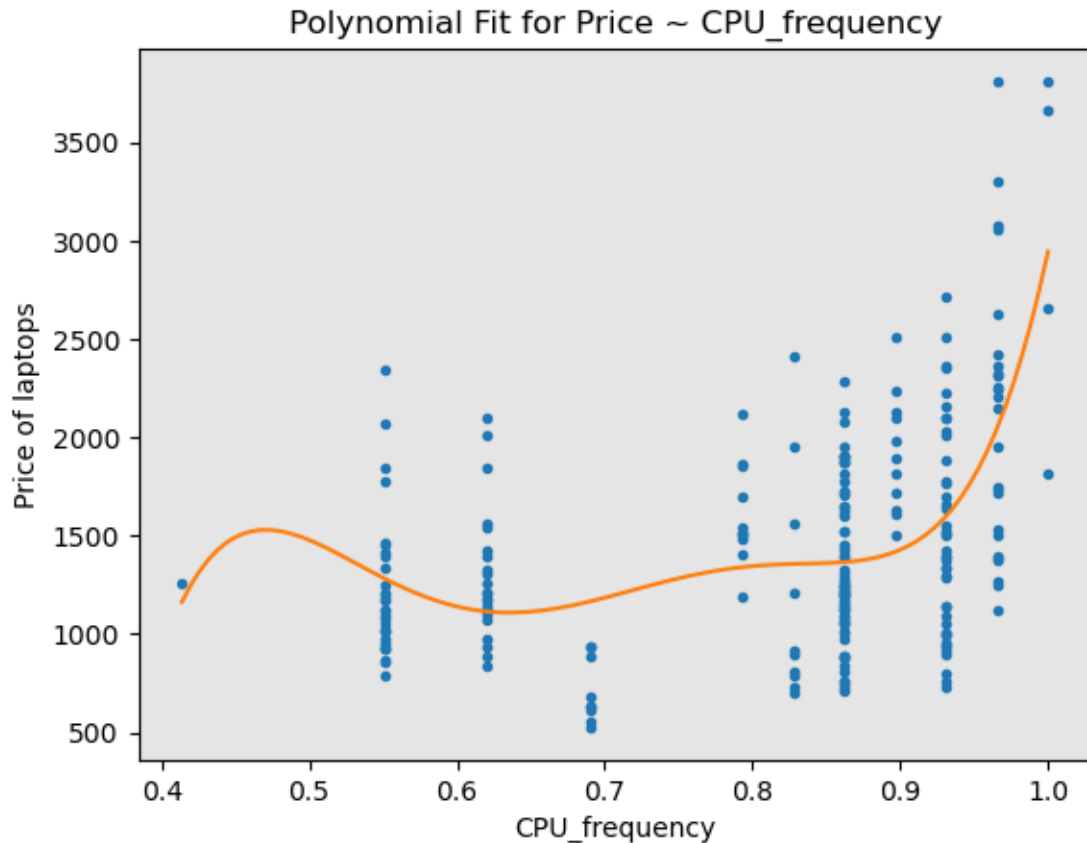
```
[27]: # Write your code below and press Shift+Enter to execute
# Call for function of degree 1
PlotPolly(p1,X,Y,'CPU_frequency')
```



```
[28]: # Write your code below and press Shift+Enter to execute  
# Call for function of degree 3  
PlotPolly(p3,X,Y,'CPU_frequency')
```



```
[29]: # Write your code below and press Shift+Enter to execute  
# Call for function of degree 5  
PlotPolly(p5,X,Y,'CPU_frequency')
```

Also, calculate the R^2 and MSE values for these fits. For polynomial functions, the function `sklearn.metrics.r2_score` will be used to calculate R^2 values.

```
[30]: # Write your code below and press Shift+Enter to execute
r_squared_1 = r2_score(Y, p1(X))
print('The R-square value for 1st degree polynomial is: ', r_squared_1)
print('The MSE value for 1st degree polynomial is: ',
      ↪mean_squared_error(Y,p1(X)))
r_squared_3 = r2_score(Y, p3(X))
print('The R-square value for 3rd degree polynomial is: ', r_squared_3)
print('The MSE value for 3rd degree polynomial is: ',
      ↪mean_squared_error(Y,p3(X)))
r_squared_5 = r2_score(Y, p5(X))
print('The R-square value for 5th degree polynomial is: ', r_squared_5)
print('The MSE value for 5th degree polynomial is: ',
      ↪mean_squared_error(Y,p5(X)))
```

```
The R-square value for 1st degree polynomial is:  0.1344436321024326
The MSE value for 1st degree polynomial is:  284583.4405868629
The R-square value for 3rd degree polynomial is:  0.2669264079653094
The MSE value for 3rd degree polynomial is:  241024.86303848823
```

The R-square value for 5th degree polynomial is: 0.3030822706442371
The MSE value for 5th degree polynomial is: 229137.29548058534

7 Task 4 - Pipeline

Create a pipeline that performs parameter scaling, Polynomial Feature generation and Linear regression. Use the set of multiple features as before to create this pipeline.

```
[35]: # Write your code below and press Shift+Enter to execute
Input=[('scale',StandardScaler()), ('polynomial',
↳PolynomialFeatures(include_bias=False)), ('model',LinearRegression())]
pipe=Pipeline(Input)
Z =
↳df[['CPU_frequency','RAM_GB','Storage_GB_SSD','CPU_core','OS','GPU','Category']]
Z = Z.astype(float)
pipe.fit(Z,Y)
ypipe=pipe.predict(Z)
```

Evaluate the MSE and R^2 values for the this predicted output.

```
[36]: # Write your code below and press Shift+Enter to execute
print('MSE for multi-variable polynomial pipeline is: ', mean_squared_error(Y,
↳ypipe))
print('R^2 for multi-variable polynomial pipeline is: ', r2_score(Y, ypipe))
```

MSE for multi-variable polynomial pipeline is: 120605.53512342437
 R^2 for multi-variable polynomial pipeline is: 0.6331800307336903

We have seen that the values of R^2 increase as we go from Single Linear Regression to Multiple Linear Regression. Further, if we go for multiple linear regression extended with polynomial features, we get an even better R^2 value.