



IBM Developer SKILLS NETWORK

Model Evaluation and Refinement

Estimated time needed: **30** minutes

Objectives

After completing this lab you will be able to:

- Evaluate and refine prediction models

Table of Contents

- [Model Evaluation \(https://#ref1\)](#)
- [Over-fitting, Under-fitting and Model Selection \(https://#ref2\)](#)
- [Ridge Regression \(https://#ref3\)](#)
- [Grid Search \(https://#ref4\)](#)

Setup

you are running the lab in your browser, so we will install the libraries using `pip`

```
In [1]: 1 #you are running the lab in your browser, so we will install the libraries
2 import piplite
3 import micropip
4 await piplite.install(['pandas'])
5 await piplite.install(['matplotlib'])
6 await piplite.install(['scipy'])
7 await piplite.install(['seaborn'])
8 await micropip.install(['ipywidgets'],keep_going=True)
9 await micropip.install(['tqdm'],keep_going=True)
```

If you run the lab locally using Anaconda, you can load the correct library and versions by uncommenting the following:

```
In [ ]: 1 #install specific version of libraries used in lab
2 #! mamba install pandas==1.3.3 -y
3 #! mamba install numpy=1.21.2 -y
4 #! mamba install sklearn=0.20.1 -y
5 #! mamba install ipywidgets=7.4.2 -y
6 #! mamba install tqdm
```

```
In [3]: 1 import pandas as pd
2 import numpy as np
```

This function will download the dataset into your browser

```
In [4]: 1 #This function will download the dataset into your browser
2
3 from pyodide.http import pyfetch
4
5 async def download(url, filename):
6     response = await pyfetch(url)
7     if response.status == 200:
8         with open(filename, "wb") as f:
9             f.write(await response.bytes())
```

```
In [5]: 1 import pandas as pd
2 import numpy as np
3
```

This dataset was hosted on IBM Cloud object. Click [HERE](https://cocl.us/DA101EN_object_storage?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01) (https://cocl.us/DA101EN_object_storage?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01) for free storage.

```
In [6]: 1 path = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/I
```

you will need to download the dataset; if you are running locally, please comment out the following

```
In [7]: 1 #you will need to download the dataset; if you are running locally, please c
2 await download(path, "auto.csv")
3 path="auto.csv"
```

```
In [8]: 1
2 df = pd.read_csv(path)
```

```
In [9]: 1 df.to_csv('module_5_auto.csv')
```

First, let's only use numeric data:

```
In [10]: 1 df=df._get_numeric_data()
2 df.head()
```

```
Out[10]:
```

	Unnamed: 0	Unnamed: 0.1	symboling	normalized- losses	wheel- base	length	width	height	curb- weight	engi s
0	0	0	3	122	88.6	0.811148	0.890278	48.8	2548	'
1	1	1	3	122	88.6	0.811148	0.890278	48.8	2548	'
2	2	2	1	122	94.5	0.822681	0.909722	52.4	2823	'
3	3	3	2	164	99.8	0.848630	0.919444	54.3	2337	'
4	4	4	2	164	99.4	0.848630	0.922222	54.3	2824	'

5 rows × 21 columns



Libraries for plotting:

```
In [11]: 1 from ipywidgets import interact, interactive, fixed, interact_manual
```

Functions for Plotting

```
In [12]: 1 def DistributionPlot(RedFunction, BlueFunction, RedName, BlueName, Title):
2     width = 12
3     height = 10
4     plt.figure(figsize=(width, height))
5
6     ax1 = sns.distplot(RedFunction, hist=False, color="r", label=RedName)
7     ax2 = sns.distplot(BlueFunction, hist=False, color="b", label=BlueName,
8
9     plt.title(Title)
10    plt.xlabel('Price (in dollars)')
11    plt.ylabel('Proportion of Cars')
12
13    plt.show()
14    plt.close()
```

```

In [13]: 1 def PollyPlot(xtrain, xtest, y_train, y_test, lr,poly_transform):
2         width = 12
3         height = 10
4         plt.figure(figsize=(width, height))
5
6
7         #training data
8         #testing data
9         # lr: linear regression object
10        #poly_transform: polynomial transformation object
11
12        xmax=max([xtrain.values.max(), xtest.values.max()])
13
14        xmin=min([xtrain.values.min(), xtest.values.min()])
15
16        x=np.arange(xmin, xmax, 0.1)
17
18
19        plt.plot(xtrain, y_train, 'ro', label='Training Data')
20        plt.plot(xtest, y_test, 'go', label='Test Data')
21        plt.plot(x, lr.predict(poly_transform.fit_transform(x.reshape(-1, 1))),
22        plt.ylim([-10000, 60000])
23        plt.ylabel('Price')
24        plt.legend()

```

Part 1: Training and Testing

An important step in testing your model is to split your data into training and testing data. We will place the target data **price** in a separate dataframe **y_data**:

```

In [14]: 1 y_data = df['price']

```

Drop price data in dataframe **x_data**:

```

In [15]: 1 x_data=df.drop('price',axis=1)

```

Now, we randomly split our data into training and testing data using the function **train_test_split**.

```

In [18]: 1 from sklearn.model_selection import train_test_split
2
3
4 x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.1)
5
6
7 print("number of test samples :", x_test.shape[0])
8 print("number of training samples:",x_train.shape[0])
9

```

```

number of test samples : 21
number of training samples: 180

```

The **test size** parameter sets the proportion of data that is split into the testing set. In the above,

the testing set is 10% of the total dataset.

Question #1):

Use the function "train_test_split" to split up the dataset such that 40% of the data samples will be utilized for testing. Set the parameter "random_state" equal to zero. The output of the function should be the following: "x_train1", "x_test1", "y_train1" and "y_test1".

```
In [19]: 1 # Write your code below and press Shift+Enter to execute
          2
          3 x_train1, x_test1, y_train1, y_test1 = train_test_split(x_data, y_data, test
          4
          5
          6 print("number of test samples :", x_test1.shape[0])
          7 print("number of training samples:",x_train1.shape[0])
          8
```

```
number of test samples : 81
number of training samples: 120
```

[Click here for the solution](#)

```
x_train1, x_test1, y_train1, y_test1 = train_test_split(x_data, y_data,
test_size=0.4, random_state=0)
print("number of test samples :", x_test1.shape[0])
print("number of training samples:",x_train1.shape[0])
```

Let's import **LinearRegression** from the module **linear_model**.

```
In [20]: 1 from sklearn.linear_model import LinearRegression
```

We create a Linear Regression object:

```
In [21]: 1 lre=LinearRegression()
```

We fit the model using the feature "horsepower":

```
In [22]: 1 lre.fit(x_train[['horsepower']], y_train)
```

```
Out[22]: LinearRegression()
```

Let's calculate the R^2 on the test data:

```
In [23]: 1 lre.score(x_test[['horsepower']], y_test)
```

```
Out[23]: 0.3635875575078824
```

We can see the R^2 is much smaller using the test data compared to the training data.

```
In [24]: 1 lre.score(x_train[['horsepower']], y_train)
```

```
Out[24]: 0.6619724197515103
```

Question #2):

Find the R^2 on the test data using 40% of the dataset for testing.

```
In [25]: 1 # Write your code below and press Shift+Enter to execute
        2 lre.fit(x_train1[['horsepower']], y_train1)
```

```
Out[25]: LinearRegression()
```

```
In [26]: 1 lre.score(x_test1[['horsepower']], y_test1)
```

```
Out[26]: 0.7139364665406973
```

```
In [27]: 1 lre.score(x_train1[['horsepower']], y_train1)
```

```
Out[27]: 0.5754067463583004
```

[Click here for the solution](#)

```
x_train1, x_test1, y_train1, y_test1 = train_test_split(x_data, y_data,
test_size=0.4, random_state=0)
lre.fit(x_train1[['horsepower']], y_train1)
lre.score(x_test1[['horsepower']], y_test1)
```

Sometimes you do not have sufficient testing data; as a result, you may want to perform cross-validation. Let's go over several methods that you can use for cross-validation.

Cross-Validation Score

Let's import **model_selection** from the module **cross_val_score**.

```
In [28]: 1 from sklearn.model_selection import cross_val_score
```

We input the object, the feature ("horsepower"), and the target data (y_data). The parameter 'cv'

determines the number of folds. In this case, it is 4.

```
In [29]: 1 Rcross = cross_val_score(lre, x_data[['horsepower']], y_data, cv=4)
```

The default scoring is R^2 . Each element in the array has the average R^2 value for the fold:

```
In [30]: 1 Rcross
```

```
Out[30]: array([0.7746232 , 0.51716687, 0.74785353, 0.04839605])
```

We can calculate the average and standard deviation of our estimate:

```
In [31]: 1 print("The mean of the folds are", Rcross.mean(), "and the standard deviatio
```

```
The mean of the folds are 0.5220099150421197 and the standard deviation is 0.29118394447560203
```

We can use negative squared error as a score by setting the parameter 'scoring' metric to 'neg_mean_squared_error'.

```
In [32]: 1 -1 * cross_val_score(lre,x_data[['horsepower']], y_data,cv=4,scoring='neg_me
```

```
Out[32]: array([20254142.84026702, 43745493.26505171, 12539630.34014929, 17561927.72247586])
```

Question #3):

Calculate the average R^2 using two folds, then find the average R^2 for the second fold utilizing the "horsepower" feature:

```
In [33]: 1 # Write your code below and press Shift+Enter to execute
2 Rcross1 = cross_val_score(lre, x_data[['horsepower']], y_data, cv=2)
3 Rcross1
```

```
Out[33]: array([0.59015621, 0.44319613])
```

```
In [34]: 1 print("The mean of the folds are", Rcross1.mean(), "and the standard deviati
```

```
The mean of the folds are 0.5166761697127429 and the standard deviation is 0.07348004195771385
```

Click here for the solution

```
Rc=cross_val_score(lre,x_data[['horsepower']], y_data,cv=2)
Rc.mean()
```

You can also use the function 'cross_val_predict' to predict the output. The function splits up the data into the specified number of folds, with one fold for testing and the other folds are used for training. First, import the function:

```
In [35]: 1 from sklearn.model_selection import cross_val_predict
```

We input the object, the feature **"horsepower"**, and the target data **y_data**. The parameter 'cv' determines the number of folds. In this case, it is 4. We can produce an output:

```
In [36]: 1 yhat = cross_val_predict(lr,x_data[['horsepower']], y_data,cv=4)
2 yhat[0:5]
```

```
Out[36]: array([14141.63807508, 14141.63807508, 20814.29423473, 12745.03562306,
14762.35027598])
```

Part 2: Overfitting, Underfitting and Model Selection

It turns out that the test data, sometimes referred to as the "out of sample data", is a much better measure of how well your model performs in the real world. One reason for this is overfitting.

Let's go over some examples. It turns out these differences are more apparent in Multiple Linear Regression and Polynomial Regression so we will explore overfitting in that context.

Let's create Multiple Linear Regression objects and train the model using **'horsepower'**, **'curb-weight'**, **'engine-size'** and **'highway-mpg'** as features.

```
In [37]: 1 lr = LinearRegression()
2 lr.fit(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']],
```

```
Out[37]: LinearRegression()
```

Prediction using training data:

```
In [38]: 1 yhat_train = lr.predict(x_train[['horsepower', 'curb-weight', 'engine-size',
2 yhat_train[0:5]
```

```
Out[38]: array([ 7426.6731551 , 28323.75090803, 14213.38819709, 4052.34146983,
34500.19124244])
```

Prediction using test data:

```
In [39]: 1 yhat_test = lr.predict(x_test[['horsepower', 'curb-weight', 'engine-size', '
2 yhat_test[0:5]
```

```
Out[39]: array([11349.35089149, 5884.11059106, 11208.6928275 , 6641.07786278,
15565.79920282])
```

Let's perform some model evaluation using our training and testing data separately. First, we

import the seaborn and matplotlib library for plotting.

```
In [40]: 1 import matplotlib.pyplot as plt
          2 %matplotlib inline
          3 import seaborn as sns
```

Let's examine the distribution of the predicted values of the training data.

In [42]:

```
1 Title = 'Distribution Plot of Predicted Value Using Training Data vs Train
2 DistributionPlot(y_train, yhat_train, "Actual Values (Train)", "Predicted Va
```

<ipython-input-12-122ce36d6117>:6: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax1 = sns.distplot(RedFunction, hist=False, color="r", label=RedName)
```

<ipython-input-12-122ce36d6117>:7: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax2 = sns.distplot(BlueFunction, hist=False, color="b", label=BlueName, ax=ax
1)
```

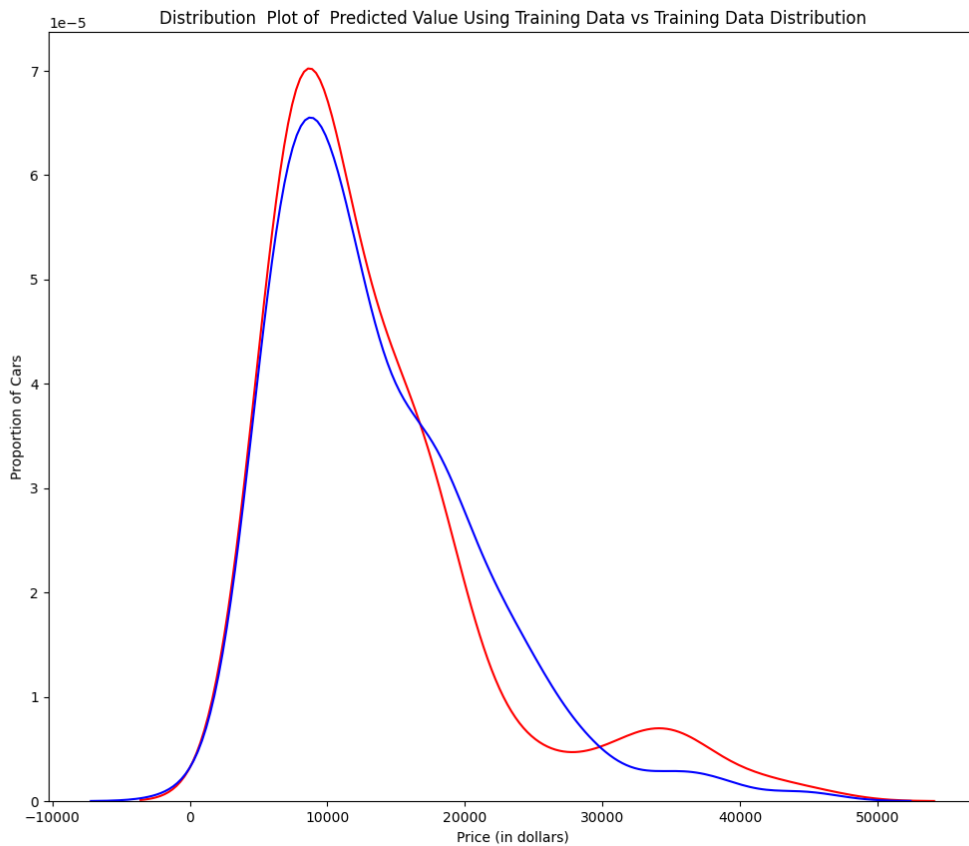


Figure 1: Plot of predicted values using the training data compared to the actual values of the training data.

So far, the model seems to be doing well in learning from the training dataset. But what happens when the model encounters new data from the testing dataset? When the model generates new values from the test data, we see the distribution of the predicted values is much different from the actual target values.

```
In [43]: 1 Title='Distribution Plot of Predicted Value Using Test Data vs Data Distri
2 DistributionPlot(y_test,yhat_test,"Actual Values (Test)","Predicted Values (
```

```
<ipython-input-12-122ce36d6117>:6: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax1 = sns.distplot(RedFunction, hist=False, color="r", label=RedName)
```

```
<ipython-input-12-122ce36d6117>:7: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax2 = sns.distplot(BlueFunction, hist=False, color="b", label=BlueName, ax=ax
1)
```

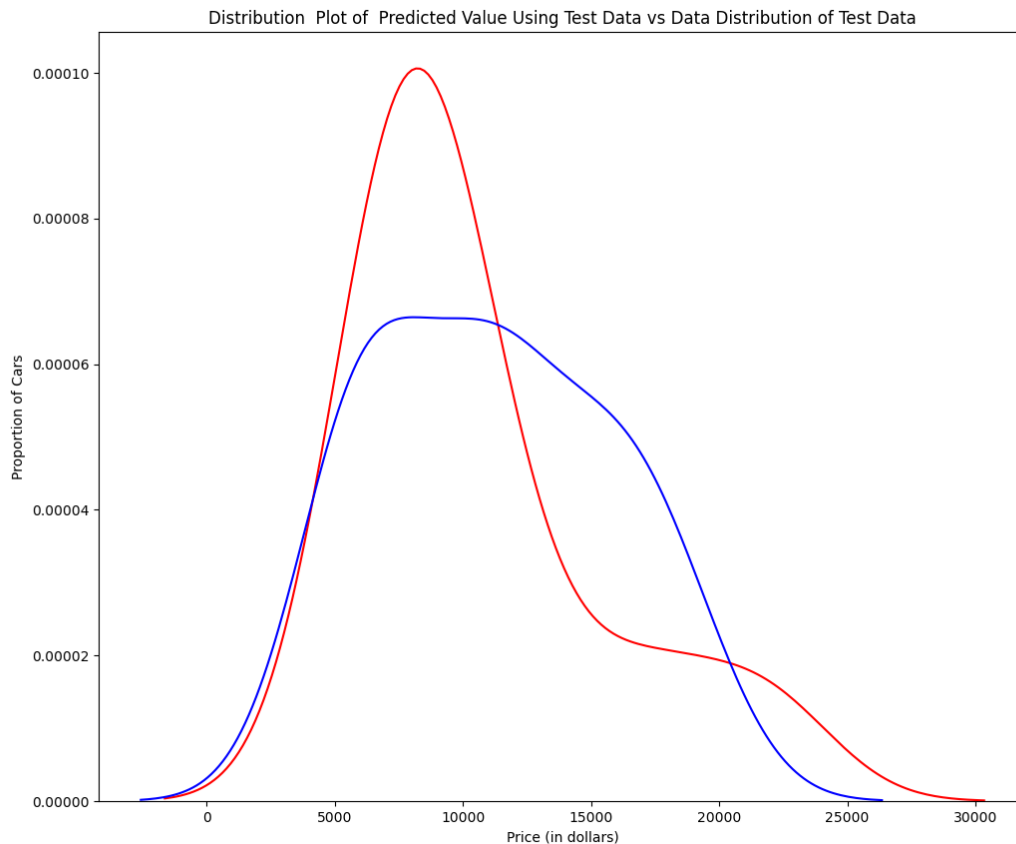


Figure 2: Plot of predicted value using the test data compared to the actual values of the test data.

Comparing Figure 1 and Figure 2, it is evident that the distribution of the test data in Figure 1 is much better at fitting the data. This difference in Figure 2 is apparent in the range of 5000 to 15,000. This is where the shape of the distribution is extremely different. Let's see if polynomial regression also exhibits a drop in the prediction accuracy when analysing the test dataset.

```
In [44]: 1 from sklearn.preprocessing import PolynomialFeatures
```

Overfitting

Overfitting occurs when the model fits the noise, but not the underlying process. Therefore, when testing your model using the test set, your model does not perform as well since it is modelling noise, not the underlying process that generated the relationship. Let's create a degree 5 polynomial model.

Let's use 55 percent of the data for training and the rest for testing:

```
In [45]: 1 x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_siz
```

We will perform a degree 5 polynomial transformation on the feature 'horsepower'.

```
In [46]: 1 pr = PolynomialFeatures(degree=5)
          2 x_train_pr = pr.fit_transform(x_train[['horsepower']])
          3 x_test_pr = pr.fit_transform(x_test[['horsepower']])
          4 pr
```

```
Out[46]: PolynomialFeatures(degree=5)
```

Now, let's create a Linear Regression model "poly" and train it.

```
In [47]: 1 poly = LinearRegression()
          2 poly.fit(x_train_pr, y_train)
```

```
Out[47]: LinearRegression()
```

We can see the output of our model using the method "predict." We assign the values to "yhat".

```
In [48]: 1 yhat = poly.predict(x_test_pr)
          2 yhat[0:5]
```

```
Out[48]: array([ 6728.58641321,  7307.91998787, 12213.73753589, 18893.37919224,
                  19996.10612156])
```

Let's take the first five predicted values and compare it to the actual targets.

```
In [49]: 1 print("Predicted values:", yhat[0:4])
          2 print("True values:", y_test[0:4].values)
```

```
Predicted values: [ 6728.58641321  7307.91998787 12213.73753589 18893.37919224]
True values: [ 6295. 10698. 13860. 13499.]
```

We will use the function "PollyPlot" that we defined at the beginning of the lab to display the training data, testing data, and the predicted function.

```
In [50]: 1 PollyPlot(x_train[['horsepower']], x_test[['horsepower']], y_train, y_test,
```

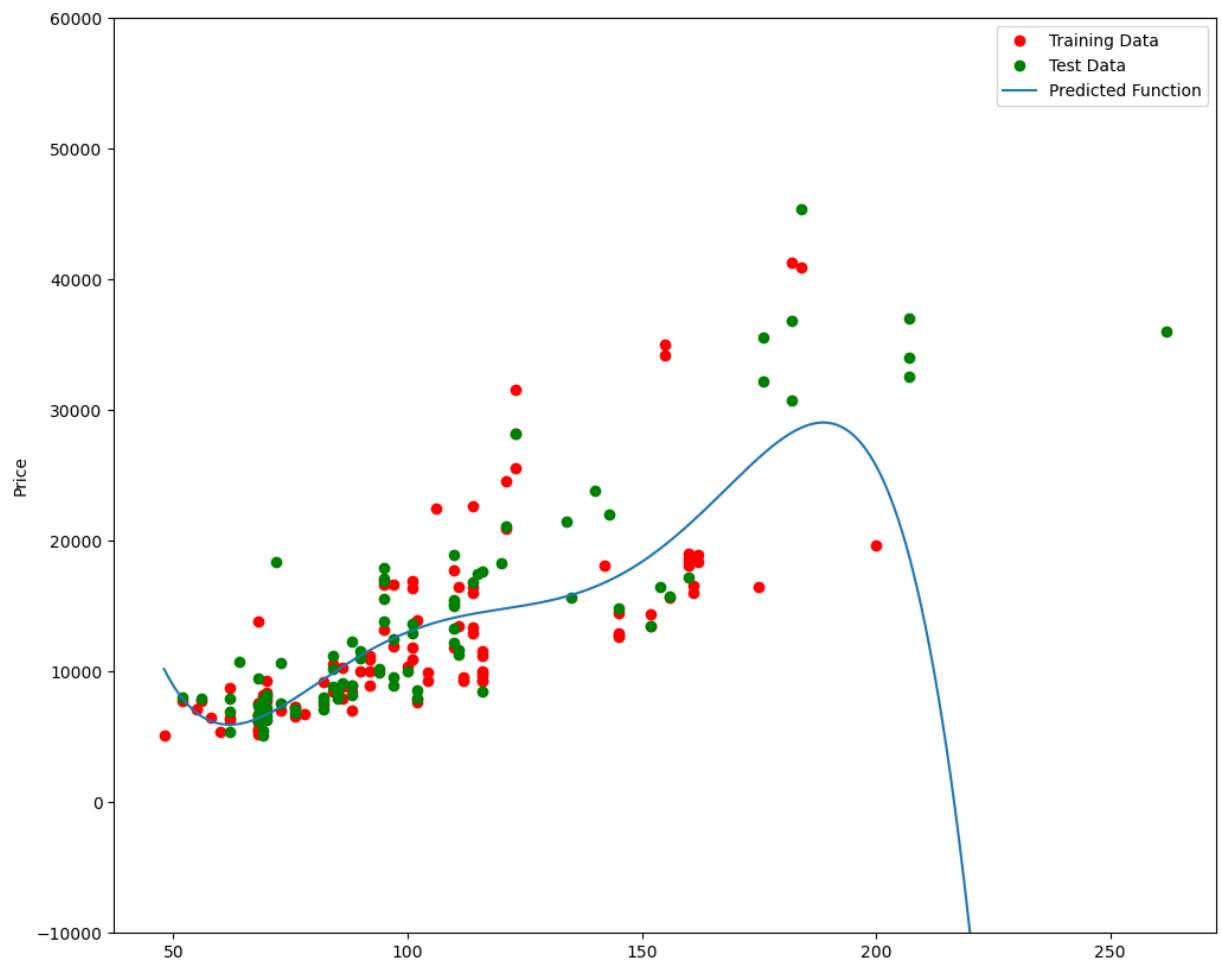


Figure 3: A polynomial regression model where red dots represent training data, green dots represent test data, and the blue line represents the model prediction.

We see that the estimated function appears to track the data but around 200 horsepower, the

function begins to diverge from the data points.

R² of the training data:

```
In [51]: 1 poly.score(x_train_pr, y_train)
```

```
Out[51]: 0.5567716897754004
```

R² of the test data:

```
In [52]: 1 poly.score(x_test_pr, y_test)
```

```
Out[52]: -29.87099623387278
```

We see the R² for the training data is 0.5567 while the R² on the test data was -29.87. The lower the R², the worse the model. A negative R² is a sign of overfitting.

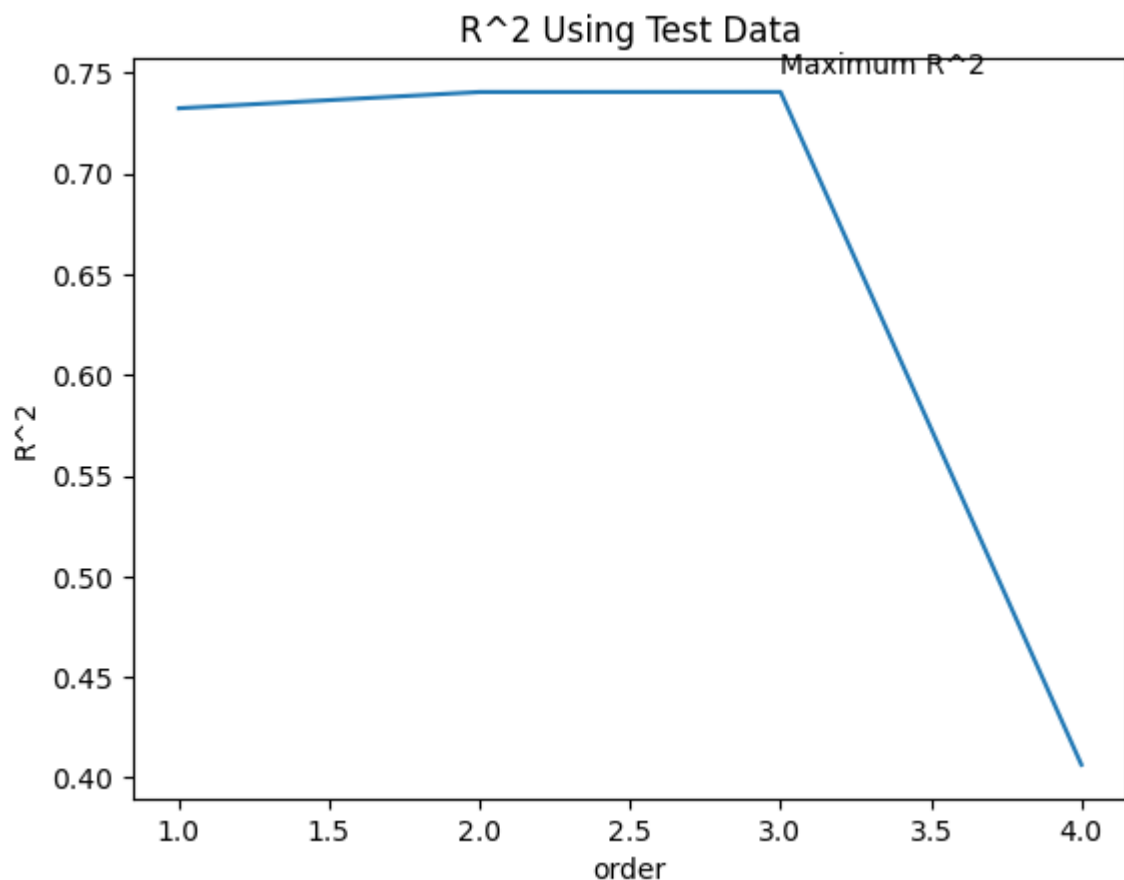
Let's see how the R² changes on the test data for different order polynomials and then plot the results:


```

In [53]: 1 Rsqu_test = []
          2
          3 order = [1, 2, 3, 4]
          4 for n in order:
          5     pr = PolynomialFeatures(degree=n)
          6
          7     x_train_pr = pr.fit_transform(x_train[['horsepower']])
          8
          9     x_test_pr = pr.fit_transform(x_test[['horsepower']])
         10
         11     lr.fit(x_train_pr, y_train)
         12
         13     Rsqu_test.append(lr.score(x_test_pr, y_test))
         14
         15 plt.plot(order, Rsqu_test)
         16 plt.xlabel('order')
         17 plt.ylabel('R^2')
         18 plt.title('R^2 Using Test Data')
         19 plt.text(3, 0.75, 'Maximum R^2 ')

```

Out[53]: Text(3, 0.75, 'Maximum R^2 ')



We see the R^2 gradually increases until an order three polynomial is used. Then, the R^2 dramatically decreases at an order four polynomial.

The following function will be used in the next section. Please run the cell below.

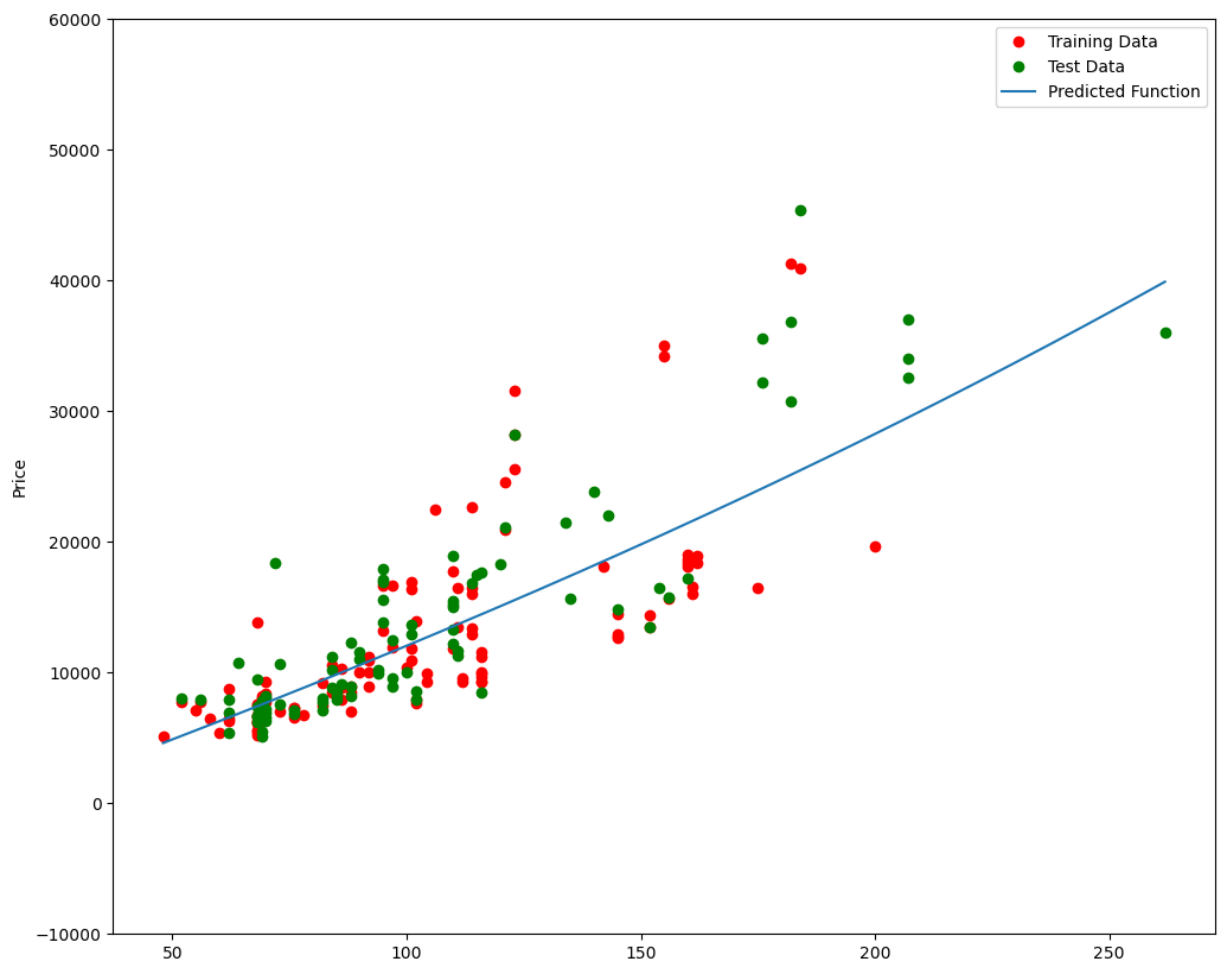
```
In [54]: 1 def f(order, test_data):
2         x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_data=test_data)
3         pr = PolynomialFeatures(degree=order)
4         x_train_pr = pr.fit_transform(x_train[['horsepower']])
5         x_test_pr = pr.fit_transform(x_test[['horsepower']])
6         poly = LinearRegression()
7         poly.fit(x_train_pr, y_train)
8         PollyPlot(x_train[['horsepower']], x_test[['horsepower']], y_train, y_test)
```

The following interface allows you to experiment with different polynomial orders and different amounts of data.

```
In [55]: 1 interact(f, order=(0, 6, 1), test_data=(0.05, 0.95, 0.05))
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

```
Out[55]: <function __main__.f(order, test_data)>
```



Question #4a):

We can perform polynomial transformations with more than one feature. Create a "PolynomialFeatures" object "pr1" of degree two.

```
In [57]: 1 # Write your code below and press Shift+Enter to execute
          2 pr1 = PolynomialFeatures(degree=2)
          3
```

[Click here for the solution](#)

```
pr1=PolynomialFeatures(degree=2)
```

Question #4b):

Transform the training and testing samples for the features 'horsepower', 'curb-weight', 'engine-size' and 'highway-mpg'. Hint: use the method "fit_transform".

```
In [59]: 1 # Write your code below and press Shift+Enter to execute
          2 x_train_pr1 = pr1.fit_transform(x_train[['horsepower', 'curb-weight', 'engine-
          3 x_test_pr1 = pr1.fit_transform(x_test[['horsepower', 'curb-weight', 'engine-si
          4 pr1
```

```
Out[59]: PolynomialFeatures()
```

[Click here for the solution](#)

```
x_train_pr1=pr1.fit_transform(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
```

```
x_test_pr1=pr1.fit_transform(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
```

Question #4c):

How many dimensions does the new feature have? Hint: use the attribute "shape".

```
In [61]: 1 # Write your code below and press Shift+Enter to execute
          2 x_train_pr1.shape
```

```
Out[61]: (110, 15)
```

[Click here for the solution](#)

```
x_train_pr1.shape #there are now 15 features
```

Question #4d):

Create a linear regression model "poly1". Train the object using the method "fit" using the polynomial features.

```
In [62]: 1 # Write your code below and press Shift+Enter to execute
          2 poly1 = LinearRegression()
          3 poly1.fit(x_train_pr1, y_train)
```

```
Out[62]: LinearRegression()
```

[Click here for the solution](#)

```
poly1=LinearRegression().fit(x_train_pr1,y_train)
```

Question #4e):

Use the method "predict" to predict an output on the polynomial features, then use the function "DistributionPlot" to display the distribution of the predicted test output vs. the actual test data.

```
In [65]: 1 # Write your code below and press Shift+Enter to execute
          2 yhat1 = poly1.predict(x_test_pr1)
          3 yhat1[0:5]
          4
```

```
Out[65]: array([ 6303.02917831, 10402.87148109, 20516.37185639, 19273.87340834,
                20555.51694564])
```

In [66]:

```
1 Title='Distribution Plot of Predicted Value Using Test Data vs Data Distri
2 DistributionPlot(y_test,yhat1,"Actual Values (Test)","Predicted Values (Test
```

<ipython-input-12-122ce36d6117>:6: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax1 = sns.distplot(RedFunction, hist=False, color="r", label=RedName)
```

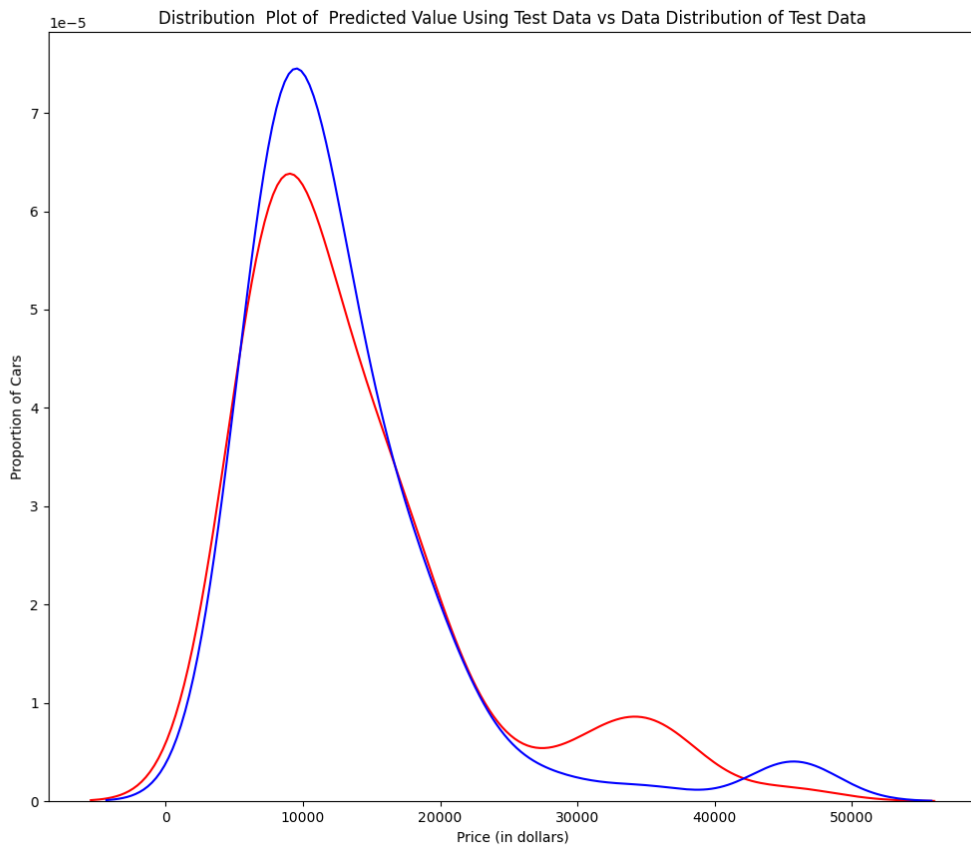
<ipython-input-12-122ce36d6117>:7: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `kdeplot` (an axes-level function for kernel density plots).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
ax2 = sns.distplot(BlueFunction, hist=False, color="b", label=BlueName, ax=ax
1)
```



[Click here for the solution](#)

```
yhat_test1=poly1.predict(x_test_pr1)
```

```
Title='Distribution Plot of Predicted Value Using Test Data vs Data Distribution of Test Data'
```

```
DistributionPlot(y_test, yhat_test1, "Actual Values (Test)", "Predicted Values (Test)", Title)
```

Question #4f):

Using the distribution plot above, describe (in words) the two regions where the predicted prices are less accurate than the actual prices.

```
In [67]: 1 # Write your code below and press Shift+Enter to execute
2 #The predicted value is higher than actual value for cars where the price $1
3 #range, conversely the predicted price is lower than the price cost in the $
4 # $40,000 range. As such the model is not as accurate in these ranges.
```

[Click here for the solution](#)

#The predicted value is higher than actual value for cars where the price is \$10,000 range, conversely the predicted price is lower than the price cost in the \$30,000 to \$40,000 range. As such the model is not as accurate in these ranges.

Part 3: Ridge Regression

In this section, we will review Ridge Regression and see how the parameter alpha changes the model. Just a note, here our test data will be used as validation data.

Let's perform a degree two polynomial transformation on our data.

```
In [68]: 1 pr=PolynomialFeatures(degree=2)
2 x_train_pr=pr.fit_transform(x_train[['horsepower', 'curb-weight', 'engine-size
3 x_test_pr=pr.fit_transform(x_test[['horsepower', 'curb-weight', 'engine-size
```

Let's import **Ridge** from the module **linear models**.

```
In [69]: 1 from sklearn.linear_model import Ridge
```

Let's create a Ridge regression object, setting the regularization parameter (alpha) to 0.1

```
In [70]: 1 RigeModel=Ridge(alpha=1)
```

Like regular regression, you can fit the model using the method **fit**.

```
In [71]: 1 RigeModel.fit(x_train_pr, y_train)
```

```
Out[71]: Ridge(alpha=1)
```

Similarly, you can obtain a prediction:

```
In [72]: 1 yhat = RigeModel.predict(x_test_pr)
```

Let's compare the first five predicted samples to our test set:

```
In [73]: 1 print('predicted:', yhat[0:4])
2 print('test set :', y_test[0:4].values)
```

```
predicted: [ 6570.82441941  9636.24891471 20949.92322738 19403.60313255]
test set : [ 6295. 10698. 13860. 13499.]
```

We select the value of alpha that minimizes the test error. To do so, we can use a for loop. We have also created a progress bar to see how many iterations we have completed so far.

```
In [74]: 1 from tqdm import tqdm
2
3 Rsqu_test = []
4 Rsqu_train = []
5 dummy1 = []
6 Alpha = 10 * np.array(range(0,1000))
7 pbar = tqdm(Alpha)
8
9 for alpha in pbar:
10     RigeModel = Ridge(alpha=alpha)
11     RigeModel.fit(x_train_pr, y_train)
12     test_score, train_score = RigeModel.score(x_test_pr, y_test), RigeModel.
13
14     pbar.set_postfix({"Test Score": test_score, "Train Score": train_score})
15
16     Rsqu_test.append(test_score)
17     Rsqu_train.append(train_score)
```

```
<ipython-input-74-e0c60797668d>:7: TqdmMonitorWarning: tqdm:disabling monitor s
upport (monitor_interval = 0) due to:
can't start new thread
  pbar = tqdm(Alpha)
100%|#####| 1000/1000 [00:02<00:00, 386.10it/s, Test Score=0.564, Train Sc
ore=0.859]
```

We can plot out the value of R^2 for different alphas:


```
In [75]: 1 width = 12
2 height = 10
3 plt.figure(figsize=(width, height))
4
5 plt.plot(Alpha, Rsqu_test, label='validation data ')
6 plt.plot(Alpha, Rsqu_train, 'r', label='training Data ')
7 plt.xlabel('alpha')
8 plt.ylabel('R^2')
9 plt.legend()
```

Out[75]: <matplotlib.legend.Legend at 0x57860a0>

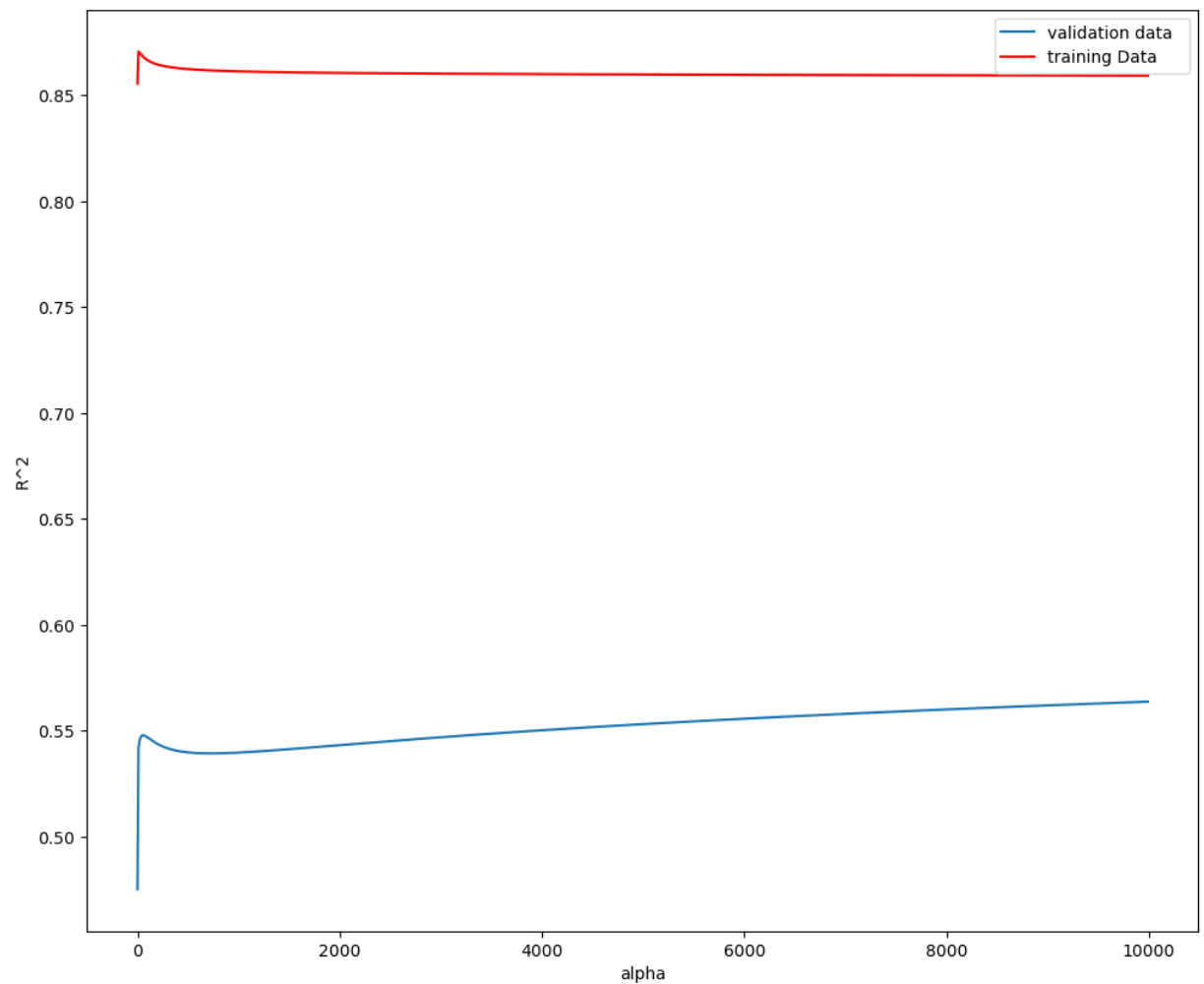


Figure 4: The blue line represents the R^2 of the validation data, and the red line represents the R^2 of the training data. The x-axis represents the different values of Alpha.

Here the model is built and tested on the same data, so the training and test data are the same.

The red line in Figure 4 represents the R^2 of the training data. As alpha increases the R^2 decreases. Therefore, as alpha increases, the model performs worse on the training data

The blue line represents the R^2 on the validation data. As the value for alpha increases, the R^2 increases and converges at a point.

Question #5):

Perform Ridge regression. Calculate the R^2 using the polynomial features, use the training data to train the model and use the test data to test the model. The parameter alpha should be set to 10.

```
In [82]: 1 RigeModel1=Ridge(alpha=10)
```

```
In [83]: 1 RigeModel1.fit(x_train_pr, y_train)
```

```
Out[83]: Ridge(alpha=10)
```

```
In [84]: 1 RigeModel.score(x_test_pr, y_test)
```

```
Out[84]: 0.5637701868993854
```

[Click here for the solution](#)

```
RigeModel = Ridge(alpha=10)
RigeModel.fit(x_train_pr, y_train)
RigeModel.score(x_test_pr, y_test)
```

Part 4: Grid Search

The term alpha is a hyperparameter. Sklearn has the class **GridSearchCV** to make the process of finding the best hyperparameter simpler.

Let's import **GridSearchCV** from the module **model_selection**.

```
In [85]: 1 from sklearn.model_selection import GridSearchCV
```

We create a dictionary of parameter values:

```
In [86]: 1 parameters1= [{'alpha': [0.001,0.1,1, 10, 100, 1000, 10000, 100000, 100000]}]  
2 parameters1
```

```
Out[86]: [{'alpha': [0.001, 0.1, 1, 10, 100, 1000, 10000, 100000, 100000]}]
```

Create a Ridge regression object:

```
In [87]: 1 RR=Ridge()  
2 RR
```

```
Out[87]: Ridge()
```

Create a ridge grid search object:

```
In [88]: 1 Grid1 = GridSearchCV(RR, parameters1,cv=4)
```

Fit the model:

```
In [89]: 1 Grid1.fit(x_data[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
```

```
Out[89]: GridSearchCV(cv=4, estimator=Ridge(),  
                    param_grid=[{'alpha': [0.001, 0.1, 1, 10, 100, 1000, 10000, 100000,  
                    100000]}])
```

The object finds the best parameter values on the validation data. We can obtain the estimator with the best parameters and assign it to the variable BestRR as follows:

```
In [90]: 1 BestRR=Grid1.best_estimator_  
2 BestRR
```

```
Out[90]: Ridge(alpha=10000)
```

We now test our model on the test data:

```
In [91]: 1 BestRR.score(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
```

```
Out[91]: 0.8411649831036152
```

Thank you for completing this lab!

Author

[Joseph Santarcangelo \(https://www.linkedin.com/in/joseph-s-50398b136/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01\)](https://www.linkedin.com/in/joseph-s-50398b136/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01)

Other Contributors

[Mahdi Noorian PhD \(https://www.linkedin.com/in/mahdi-noorian-58219234/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01\)](https://www.linkedin.com/in/mahdi-noorian-58219234/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01)

Bahare Talayian

Eric Xiao

Steven Dong

Parizad

Hima Vasudevan

[Fiorella Wenver \(https://www.linkedin.com/in/fiorellawever/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01\)](https://www.linkedin.com/in/fiorellawever/?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=1000655SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01)

[Yi Yao \(https://www.linkedin.com/in/yi-leng-yao-84451275/\)](https://www.linkedin.com/in/yi-leng-yao-84451275/)

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2022-08-23	2.4	Malika	Updated packages by adding keep_going=True
2020-10-30	2.3	Lakshmi	Changed URL of csv
2020-10-05	2.2	Lakshmi	Removed unused library imports
2020-09-14	2.1	Lakshmi	Made changes in OverFitting section
2020-08-27	2.0	Lavanya	Moved lab to course repo in GitLab

© IBM Corporation 2020. All rights reserved.

