

# B.Sc. (Hons.) Computer Science

## Data Analysis and Visualization Project

Name : Khushal Rastogi

Roll No. : 19013570016 (2k19/CS/60)

Year & Sem : 3<sup>rd</sup> Year & 5<sup>th</sup> Sem

### >>Problem Statement:

Project is based on the problem of “**Predicting prices of used cars?**”.

#### *Real Life Scenario:*

- Assume you have a friend Sam, who wants to sell his car.
- But the problem is, he doesn't know how much he should sell his car for?
- He wants to sell it for as much as he can; but he also wants to set the price reasonably so someone would buy it.
- The Big Question is “How can we help him determine the best price for his car?”
- Let us try to clearly define some of his problems.
  - *Is there any data on the prices of some cars & their characteristics?*
  - *What features of cars affect their prices?*  
*Color? Brand? Horsepower? Something Else...?*
  - *How much can a particular feature affect the price relative to others?*
- These are some of the Questions we can start thinking about at the start!

### >>Understanding the Data:

- The Dataset used is an open dataset, by Jeffrey C. Schlemmer.

- This dataset is in CSV (Comma Separated Values) format, where each of the values are separated by commas.
- There are a total of **201 different entries** (rows) in the dataset.
- There are **26 different features** (columns) in the dataset.
  - Some of them are numerical variables like wheel-base, length, engine-size, horsepower, price, etc.
  - While some are categorical variables like fuel-type, body-style, brand, engine location, etc.
- **Description of Attributes** is given in the table below:

No.	Attribute name	attribute range	No.	Attribute name	attribute range
1	symboling	-3, -2, -1, 0, 1, 2, 3.	14	curb-weight	continuous from 1488 to 4066.
2	normalized-losses	continuous from 65 to 256.	15	engine-type	dohc, dohcvt, l, ohc, ohcf, ohcv, rotor.
3	make	audi, bmw, etc.	16	num-of-cylinders	eight, five, four, six, three, twelve, two.
4	fuel-type	diesel, gas.	17	engine-size	continuous from 61 to 326.
5	aspiration	std, turbo.	18	fuel-system	1bbl, 2bbl, 4bbl, idi, mfi, mpfi, spdi, spfi.
6	num-of-doors	four, two.	19	bore	continuous from 2.54 to 3.94.
7	body-style	hardtop, wagon, etc.	20	stroke	continuous from 2.07 to 4.17.
8	drive-wheels	4wd, fwd, rwd.	21	compression-ratio	continuous from 7 to 23.
9	engine-location	front, rear.	22	horsepower	continuous from 48 to 288.
10	wheel-base	continuous from 86.6 to 120.9.	23	peak-rpm	continuous from 4150 to 6600.
11	length	continuous from 141.1 to 208.1.	24	city-mpg	continuous from 13 to 49.
12	width	continuous from 60.3 to 72.3.	25	highway-mpg	continuous from 16 to 54.
13	height	continuous from 47.8 to 59.8.	26	price	continuous from 5118 to 45400.



# Data Analysis and Visualization Project

**Khusal Rastogi**  
**B.Sc. (Hons.) CS**  
**3rd Year, 5th Sem**

**Topic :- Predicting Prices of Used Cars !!!**

## Table of Contents

1. Data Acquisition
2. Data Wrangling
3. Exploratory Data Analysis
4. Model Development
5. Model Evaluation and Refinement
6. Conclusion

## Data Acquisition

There are various formats for a dataset: csv, json, xls etc. The dataset can be stored in different places, on your local machine or sometimes online.  
In our case, the Automobile Dataset is an online source, and it is in a CSV (comma separated value) format.

- Data source:- <https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data>
- Data type:- csv file

The **Pandas** Library is a useful tool that enables us to read various datasets into a dataframe. Also, The **Numpy** Library can be used to perform a wide variety of mathematical operations on arrays.

```
In [1]: # Import pandas library
import pandas as pd
import numpy as np
```

## Read Data

We use `pandas.read_csv()` function to read the csv file. In the brackets, we put the file path along with a quotation mark so that pandas will read the file into a dataframe from that address.  
Because the data does not include headers, we can add an argument `headers = None` inside the `read_csv()` method so that pandas will not automatically set the first row as a header.

```
In [2]: # Reading the online file mentioned in the URL above, and assigning it to variable "df".
path = "C://Users/Akash/Desktop/Cars Data Project/Raw_Automobile.csv"
df = pd.read_csv(path, header=None)
```

After reading the dataset, we can use the `dataframe.head(n)` method to check the top `n` rows of the dataframe, where `n` is an integer. Contrary to `dataframe.head(n)`, `dataframe.tail(n)` will show you the bottom `n` rows of the dataframe.

```
In [3]: # Showing the first 5 rows using dataframe.head() method.
print("The first 5 rows of the dataframe:\n")
df.head(5)
```

The first 5 rows of the dataframe.

	0	1	2	3	4	5	6	7	8	9	...	16	17	18	19	20	21	22	23	24	25
0	0	1	2	3	4	5	6	7	8	9	...	16	17	18	19	20	21	22	23	24	25
1	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21	27	13495
2	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21	27	16500
3	1	?	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154	5000	19	26	16500

5 rows x 26 columns

```
In [4]: # Showing the last 5 rows using dataframe.tail() method.
print("The last 5 rows of the dataframe:\n")
df.tail(5)
```

The last 5 rows of the dataframe.

	0	1	2	3	4	5	6	7	8	9	...	16	17	18	19	20	21	22	23	24	25
201	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	9.5	114	5400	23	28	16845
202	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	8.7	160	5300	19	25	19045
203	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	173	mpfi	3.58	2.87	8.8	134	5500	18	23	21485
204	-1	95	volvo	diesel	turbo	four	sedan	rwd	front	109.1	...	145	idi	3.01	3.40	23.0	106	4800	26	27	22470
205	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	9.5	114	5400	19	25	22625

5 rows x 26 columns

## Adding Headers

Take a look at the dataset. Pandas automatically set the header with an integer starting from 0.  
To better describe our data, we can introduce a header. This information is available at: <https://archive.ics.uci.edu/ml/datasets/Automobile>.

Thus, we have to add headers manually.

First, we create a list "headers" that include all column names in order. Then, we use `dataframe.columns = headers` to replace the headers with the list we created.

```
In [5]: # Created headers list
headers = ["symboling", "normalized-losses", "make", "fuel-type", "aspiration", "num-of-doors", "body-style", "drive-wheels", "engine-location", "wheel-base", "length", "width", "height", "curb-weight", "engine-type", "num-of-cylinders", "fuel-system", "bore", "stroke", "compression-ratio", "horsepower", "peak-rpm", "city-mpg", "highway-mpg", "price"]
print(headers, headers)
```

We replace headers with "headers":

```
In [6]: df.columns = headers
#Removing the previous header.
df.drop(index=0, inplace=True)
df.head(10)
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore	stroke	compression-ratio
1	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
2	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
3	1	?	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	
4	2	164	audi	gas	std	four	sedan	4wd	front	99.8	...	109	mpfi	3.19	3.40	
5	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	
6	2	?	audi	gas	std	two	sedan	4wd	front	99.8	...	136	mpfi	3.19	3.40	
7	1	158	audi	gas	std	four	sedan	4wd	front	105.8	...	136	mpfi	3.19	3.40	
8	1	?	audi	gas	std	four	wagon	4wd	front	105.8	...	136	mpfi	3.19	3.40	
9	1	158	audi	gas	turbo	four	sedan	4wd	front	105.8	...	131	mpfi	3.13	3.40	
10	0	?	audi	gas	turbo	two	hatchback	4wd	front	99.5	...	131	mpfi	3.13	3.40	

10 rows x 26 columns

Now, we have successfully read the raw dataset and added the correct headers into the dataframe.

## Data Wrangling

### Table of Contents

1. Identify and handle missing values
• Identify missing values
• Deal with missing values
• Correct data format
2. Basic Insight of Dataset
3. Data normalization
4. Binning
5. Indicator variable

Data wrangling is the process of converting data from the initial format to a format that may be better for analysis.

As we can see, several question marks appeared in the dataframe; these are missing values which may hinder our further analysis.

So, how do we identify all those missing values and deal with them?

**How to work with missing data?**

Steps for working with missing data:

1. Identify missing data
2. Deal with missing data
3. Correct data format

```
In [7]: df.head()
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore	stroke	compression-ratio
1	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
2	3	?	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
3	1	?	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	
4	2	164	audi	gas	std	four	sedan	4wd	front	99.8	...	109	mpfi	3.19	3.40	
5	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	

5 rows x 26 columns

## Identify and handle missing values

### Identify missing values

Convert "?" to NaN

In the `cars` dataset, missing data comes with the question mark "?". We replace "?" with `NaN` (Not a Number), Python's default missing value marker for reasons of computational speed and convenience. Here we use the function:

```
df.replace(A, B, inplace = True)
# replace "A" by B
df.replace("?", np.nan, inplace = True)
df.isnull().sum()
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore	stroke	compression-ratio
1	3	NaN	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
2	3	NaN	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	
3	1	NaN	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	
4	2	164	audi	gas	std	four	sedan	4wd	front	99.8	...	109	mpfi	3.19	3.40	
5	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	

5 rows x 26 columns

### Evaluating for Missing Data

The output is a boolean value indicating whether the value that is passed into the argument is in fact missing data.

```
In [9]: missing_data = df.isnull()
missing_data.head(5)
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore	stroke	compression-ratio
1	False	True	False	False	False	False	False	False	False	...	False	False	False	False	False	False
2	False	True	False	False	False	False	False	False	False	...	False	False	False	False	False	False
3	False	True	False	False	False	False	False	False	False	...	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
5	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False

5 rows x 26 columns

"True" means the value is a missing value while "False" means the value is not a missing value.

### Count missing values in each column

In the body of the for loop the method `value_counts()` counts the number of "True" values.

```
In [10]: for column in missing_data.columns.values.tolist():
print(">>>column")
print(missing_data[column].value_counts())
print("=====")

>> symboling
False    205
Name: symboling, dtype: int64
-----
>> normalized-losses
False    164
True      41
Name: normalized-losses, dtype: int64
-----
>> make
False    205
Name: make, dtype: int64
-----
>> fuel-type
False    203
Name: fuel-type, dtype: int64
-----
>> aspiration
False    205
Name: aspiration, dtype: int64
-----
>> num-of-doors
False    203
True       2
Name: num-of-doors, dtype: int64
-----
>> body-style
False    205
Name: body-style, dtype: int64
-----
>> drive-wheels
False    205
Name: drive-wheels, dtype: int64
-----
>> engine-location
False    205
Name: engine-location, dtype: int64
-----
>> wheel-base
False    205
Name: wheel-base, dtype: int64
-----
>> length
False    205
Name: length, dtype: int64
-----
>> width
False    205
Name: width, dtype: int64
-----
>> height
False    205
Name: height, dtype: int64
-----
>> curb-weight
False    205
Name: curb-weight, dtype: int64
-----
>> engine-type
False    203
Name: engine-type, dtype: int64
-----
>> num-of-cylinders
False    201
Name: num-of-cylinders, dtype: int64
-----
>> fuel-system
False    201
Name: fuel-system, dtype: int64
-----
>> bore
False    201
Name: bore, dtype: int64
-----
>> stroke
True      4
Name: stroke, dtype: int64
-----
>> compression-ratio
False    205
Name: compression-ratio, dtype: int64
-----
>> horsepower
False    203
True       2
Name: horsepower, dtype: int64
-----
>> peak-rpm
False    203
True       2
Name: peak-rpm, dtype: int64
-----
>> city-mpg
False    205
Name: city-mpg, dtype: int64
-----
>> highway-mpg
False    205
Name: highway-mpg, dtype: int64
-----
>> price
False    201
True       4
Name: price, dtype: int64
-----
```

Based on the summary above, each column has 205 rows of data and seven of the columns containing missing data:

1. "normalized-losses": 41 missing data
2. "num-of-doors": 2 missing data
3. "bore": 4 missing data
4. "stroke": 4 missing data
5. "horsepower": 2 missing data
6. "peak-rpm": 2 missing data
7. "price": 4 missing data

### Deal with missing data

**How to deal with missing data?**

- a. Drop the whole row
  - b. Drop the whole column
- a. Replace data
  - b. Replace it by mean
  - c. Replace it by frequency
  - c. Replace it based on other functions

Whole columns should be dropped only if most entries in the column are empty. In our dataset, none of the columns are empty enough to drop entirely. We have some freedom in choosing which method to replace data; however, some methods may seem more reasonable than others. We will apply each method to many different columns.

**Replace by mean:**

- "normalized-losses": 41 missing data, replace them with mean
- "stroke": 4 missing data, replace them with mean
- "bore": 4 missing data, replace them with mean
- "horsepower": 2 missing data, replace them with mean
- "peak-rpm": 2 missing data, replace them with mean

**Replace by frequency:**

- "num-of-doors": 2 missing data, replace them with "four".
  - Reason: 84% sedans are four doors. Since four doors is most frequent, it is most likely to occur.

**Drop the whole row:**

- "price": 4 missing data, simply delete the whole row
  - Reason: price is what we want to predict. Any data entry without price data cannot be used for prediction; therefore any row now without price data is not useful to us.

**Calculate the mean value for the "normalized-losses" column**

```
In [11]: avg_norm_loss = df["normalized-losses"].astype("float").mean(axis=0)
print("Average of normalized-losses:", avg_norm_loss)
```

Average of normalized-losses: 122.0

**Replace "NaN" with mean value in "normalized-losses" column**

```
In [12]: df["normalized-losses"].replace(np.nan, avg_norm_loss, inplace=True)
```

**Calculate the mean value for the "bore" column**

```
In [13]: avg_bore=df["bore"].astype('float').mean(axis=0)
print("Average of bore:", avg_bore)
```

Average of bore: 3.3297512437819957

**Replace "NaN" with the mean value in the "bore" column**

```
In [14]: df["bore"].replace(np.nan, avg_bore, inplace=True)
```

**Calculate the mean value for the "stroke" column**

```
In [15]: avg_stroke = df["stroke"].astype("float").mean(axis = 0)
print("Average of stroke:", avg_stroke)
```

Average of stroke: 3.2554228855721337

**Replace "NaN" with the mean value in the "stroke" column**

```
In [16]: df["stroke"].replace(np.nan, avg_stroke, inplace = True)
```

**Calculate the mean value for the "horsepower" column**

```
In [17]: avg_horsepower = df["horsepower"].astype("float").mean(axis=0)
print("Average horsepower:", avg_horsepower)
```

Average horsepower: 104.25615763546799

**Replace "NaN" with the mean value in the "horsepower" column**

```
In [18]: df["horsepower"].replace(np.nan, avg_horsepower, inplace=True)
```

**Calculate the mean value for "peak-rpm" column**

```
In [19]: avg_peakrpm=df["peak-rpm"].astype("float").mean(axis=0)
print("Average peak-rpm:", avg_peakrpm)
```

Average peak rpm: 5125.369459128679

**Replace "NaN" with the mean value in the "peak-rpm" column**

```
In [20]: df["peak-rpm"].replace(np.nan, avg_peakrpm, inplace=True)
```

To see which values are present in a particular column, we can use the `value_counts()` method:

```
In [21]: df["num-of-doors"].value_counts()
```

	four	two
1	114	89
2	114	89
3	114	89
4	114	89
5	114	89

We can see that four doors are the most common type.

We can also use the `idxmax()` method to calculate the most common type automatically:

```
In [22]: df["num-of-doors"].value_counts().idxmax()
```

"four"

The replacement procedure is very similar to what we have seen previously.

```
In [23]: #replace the missing "num-of-doors" values by the most frequent
df["num-of-doors"].replace(np.nan, "four", inplace=True)
```

Finally, let's drop all rows that do not have price data:

```
In [24]: # simply drop whole row with NaN in "price" column
df.dropna(subset=["price"], axis=0, inplace=True)
# reset index, because we dropped two rows
df.reset_index(drop=True, inplace=True)
```

```
In [25]: df.head()
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	...	engine-size	fuel-system	bore
--	-----------	-------------------	------	-----------	------------	--------------	------------	--------------	-----------------	------------	-----	-------------	-------------	------



```
0 0.811148 0.890278 0.816054
2 0.822681 0.909722 0.876254
3 0.848630 0.919444 0.908027
4 0.848630 0.922222 0.908027
```

Here we can see we've normalized 'length', 'width' and 'height' in the range of [0,1].

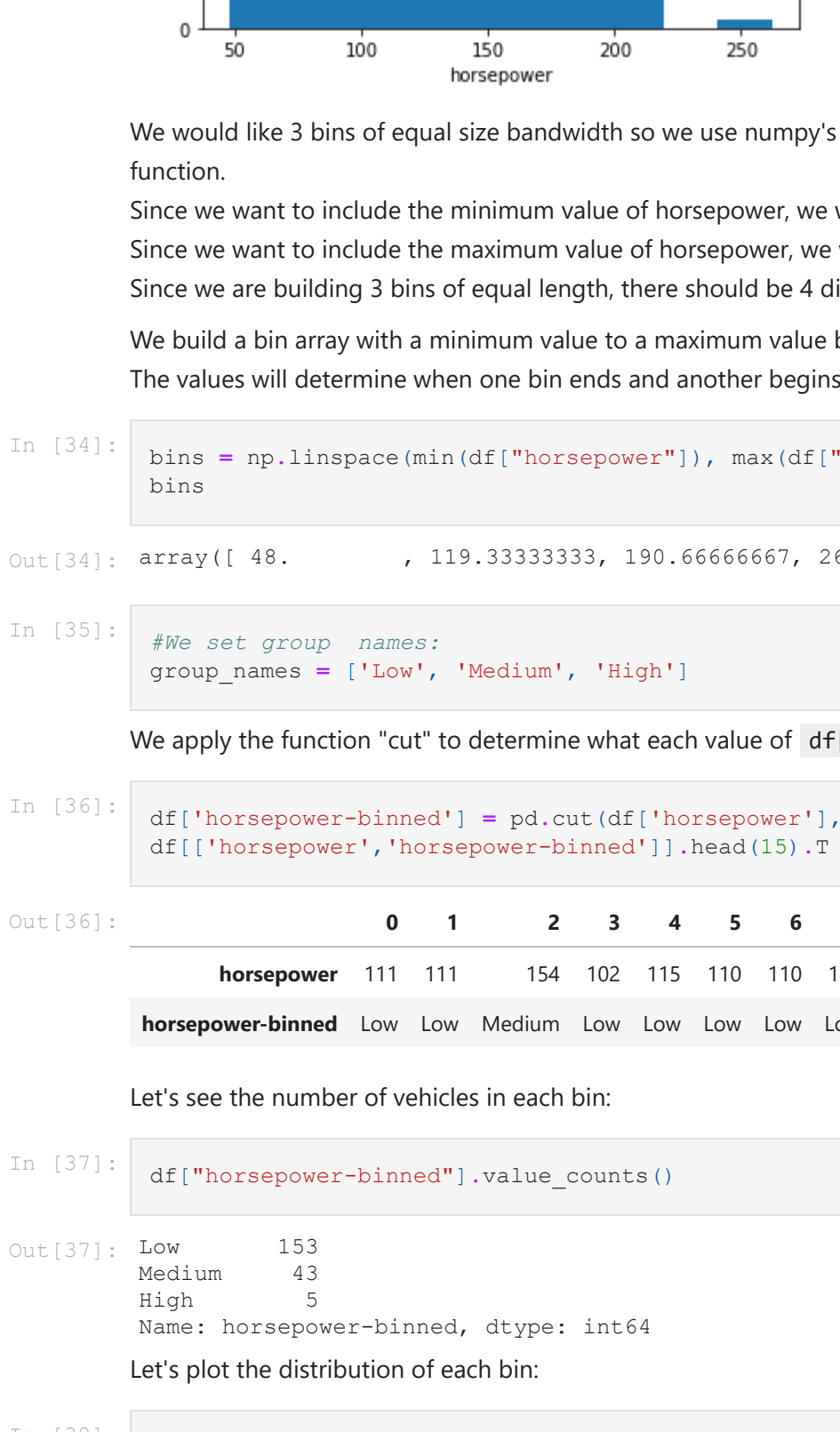
## Binning

### Why binning?

Binning is a process of transforming continuous numerical variables into discrete categorical 'bins' for grouped analysis.

In our dataset, 'horsepower' is a real valued variable ranging from 48 to 288 and it has 57 unique values. What if we only care about the price difference between cars with high horsepower, medium horsepower, and little horsepower (3 types)? Can we rearrange them into three 'bins' to simplify analysis?

Let's Plot the Histogram of horsepower to see what the distribution of horsepower looks like.



We would like 3 bins of equal size bandwidth so we use numpy's `linspace(start_value, end_value, numbers_generated)` function.

Since we want to include the minimum value of horsepower, we want to set `start_value = min(df['horsepower'])`.

Since we want to include the maximum value of horsepower, we want to set `end_value = max(df['horsepower'])`.

Since we are building 3 bins of equal length, there should be 4 dividers, so `numbers_generated = 4`.

We build a bin array with a minimum value to a maximum value by using the bandwidth calculated above

The values will determine when one bin ends and another begins.

```
In [34]: bins = np.linspace(min(df["horsepower"]), max(df["horsepower"]), 4)
bins
array([ 49.          , 119.33333333, 190.66666667, 262.          ])
```

```
In [35]: #We set group names:
group_names = ['Low', 'Medium', 'High']
```

We apply the function 'cut' to determine what each value of `df['horsepower']` belongs to.

```
In [36]: df["horsepower-binned"] = pd.cut(df["horsepower"], bins, labels=group_names, include_lowest=True)
df[["horsepower", "horsepower-binned"]].head(15).T
```

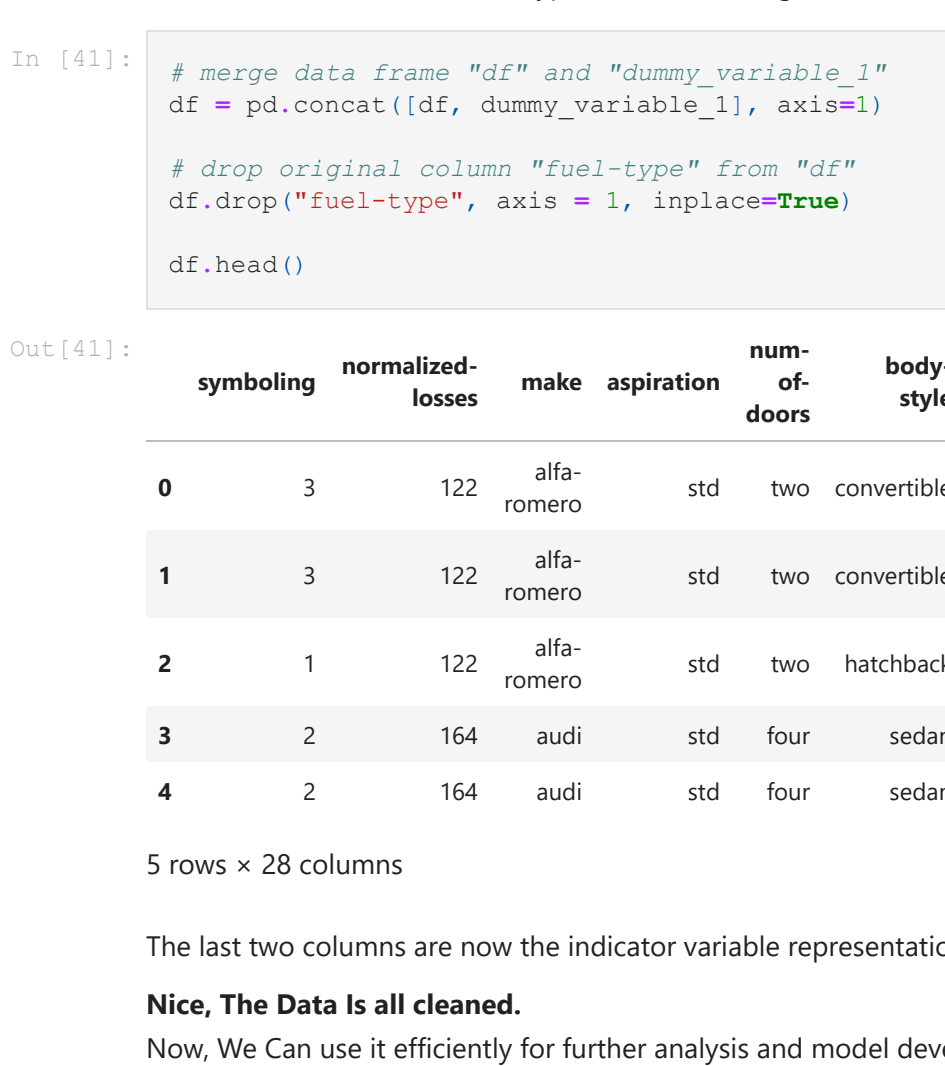
horsepower	111	1	154	102	115	110	110	7	8	10	101	11	12	13	182
horsepower-binned	Low	Low	Medium	Low	Low	Low	Low	Medium	Low	Low	Medium	Medium	Medium	Medium	Medium

Let's see the number of vehicles in each bin:

```
In [37]: df["horsepower-binned"].value_counts()

Low      153
Medium   43
High      8
Name: horsepower-binned, dtype: int64
```

Let's plot the distribution of each bin:



We successfully narrowed down the intervals from 57 to 3!

## Indicator Variable (or Dummy Variable)

An indicator variable (or dummy variable) is a numerical variable used to label categories. They are called 'dummies' because the numbers themselves don't have inherent meaning.

We use indicator variables so we can use categorical variables for regression analysis.

### Creating an indicator variable for "fuel-type".

We see the column "fuel-type" has two unique values, "gas" or "diesel". Regression doesn't understand words, only numbers.

To use this attribute in regression analysis, we convert "fuel-type" into indicator variables.

We will use pandas' method 'get\_dummies' to assign numerical values to different categories of fuel type.

```
In [39]: dummy_variable_1 = pd.get_dummies(df["fuel-type"])
dummy_variable_1.head()
```

	diesel	gas
0	0	1
1	0	1
2	0	1
3	0	1
4	0	1

Change the column names for clarity:

```
In [40]: #Changing the column names for clarity.
dummy_variable_1.rename(columns={"gas":'fuel-type=gas', 'diesel':'fuel-type=diesel'}, inplace=True)
dummy_variable_1.head()
```

	fuel-type=diesel	fuel-type=gas
0	0	1
1	0	1
2	0	1
3	0	1
4	0	1

In the dataframe, column 'fuel-type' has values for 'gas' and 'diesel' as 0s and 1s now.

```
In [41]: # merge data frame "df" and "dummy_variable_1"
df = pd.concat([df, dummy_variable_1], axis=1)

# drop original column "fuel-type" from "df"
df.drop("fuel-type", axis = 1, inplace=True)
df.head()
```

	symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	...	stroke	compression-ratio	horsepower
0	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	2.68	9.0	111
1	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	2.68	9.0	111
2	1	122	alfa-romero	std	two	hatchback	rwd	front	94.5	0.822681	...	3.47	9.0	154
3	2	164	audi	std	four	sedan	fwd	front	99.8	0.848630	...	3.40	10.0	102
4	2	164	audi	std	four	sedan	4wd	front	99.4	0.848630	...	3.40	8.0	115

5 rows x 28 columns

The last two columns are now the indicator variable representation of the fuel-type variable. They're all 0s and 1s now.

**Now, The Data is all cleaned.**

**Nice, We can use it efficiently for further analysis and model development.**

## Exploratory Data Analysis

### Table of Contents

1. Importing Cleaned Data
2. Analyzing Individual Feature Patterns using Visualization
3. Descriptive Statistical Analysis
4. Basics of Grouping
5. Correlation
6. ANOVA

## What are the main characteristics that have the most impact on the car price?

### Cleaned Data

```
In [42]: df.head()
```

	symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	...	stroke	compression-ratio	horsepower
0	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	2.68	9.0	111
1	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	...	2.68	9.0	111
2	1	122	alfa-romero	std	two	hatchback	rwd	front	94.5	0.822681	...	3.47	9.0	154
3	2	164	audi	std	four	sedan	fwd	front	99.8	0.848630	...	3.40	10.0	102
4	2	164	audi	std	four	sedan	4wd	front	99.4	0.848630	...	3.40	8.0	115

5 rows x 28 columns

### Analyzing Individual Feature Patterns Using Visualization

Importing visualization packages "Matplotlib" and "Seaborn".

```
In [43]: import matplotlib.pyplot as plt
import seaborn as sns

# list the data types for each column
print(df.dtypes)
```

symboling	int64
normalized-losses	int32
make	object
aspiration	object
num-of-doors	object
body-style	object
drive-wheels	object
engine-location	object
wheel-base	float64
length	float64
width	float64
height	float64
curb-weight	int64
engine-type	object
num-of-cylinders	object
engine-size	int64
fuel-system	object
bore	float64
stroke	float64
compression-ratio	float64
horsepower	int64
peak-rpm	float64
city-mpg	int64
highway-mpg	int64
price	float64
horsepower-binned	category
fuel-type=diesel	uint8
fuel-type=gas	uint8
dtype:	object

We can calculate the correlation between variables of type 'int64' or 'float64' using the method "corr".

```
In [45]: df.corr()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.140019	-0.008153	-0.182196	
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.068602	-0.373737	0.099404	0.112360	-0.029862	0.055045	-0.147113	
wheel-base	-0.535987	-0.056661	1.000000	0.876024	0.086802	-0.373737	0.099404	0.112360	-0.029862	0.055045	-0.147113	
length	-0.365404	0.019424	0.876024	1.000000	0.857170	0.490742	0.782097	0.572027	0.493244	0.158018	0.250313	
width	-0.242423	0.068602	0.814507	0.857170	1.000000	0.306002	0.866201	0.729436	0.544885	0.188822	0.189867	
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.180449	-0.060663	0.259737	
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.644060	0.167438	0.156433	
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.572609	0.205928	0.002889	
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.000000	-0.055390	0.001263	
stroke	-0.008153	0.055045	0.158018	0.123952	0.188822	-0.060663	0.167438	0.205928	-0.055390	1.000000	0.187871	
compression-ratio	-0.182196	-0.147113	0.250313	0.159733	0.189867	-0.259737	0.156433	0.028889	0.001263	0.187871	1.000000	
horsepower	0.078160	-0.127130	0.371178	0.579795	0.615056	-0.087001	0.757981	0.822668	0.001569	0.098128	-0.214489	
peak-rpm	0.279740	0.239543	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.267392	-0.063561	-0.435780	
city-mpg	-0.035257	-0.225016	-0.470006	-0.665192	-0.633531	-0.049800	-0.749543	-0.650546	-0.582027	-0.033956	0.331425	
highway-mpg	0.038233	-0.181877	-0.543304	-0.698142	-0.690635	-0.104812	-0.794889	-0.679571	-0.591309	-0.04636	0.268465	
price	-0.082391	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.543155	0.082269	0.077107	
fuel-type=diesel	-0.196735	-0.101546	0.307237	0.211187	0.244356	0.281578	0.221046	0.070779	0.054458	0.241064	0.985321	
fuel-type=gas	0.196735	0.101546	-0.307237	-0.211187	-0.244356	-0.281578	-0.221046	-0.070779	-0.054458	-0.241064	-0.985321	

The diagonal elements are always one.

```
In [46]: #Find the correlation between the following columns: bore, stroke, compression-ratio, and horsepower.
df[["bore", "stroke", "compression-ratio", "horsepower"]].corr()
```

	bore	stroke	compression-ratio	horsepower
bore	1.000000	-0.055390	0.001263	0.566903
stroke	-0.055390	1.000000	0.187871	0.098128
compression-ratio	0.001263	0.187871	1.000000	-0.214489
horsepower	0.566903	0.098128	-0.214489	1.000000

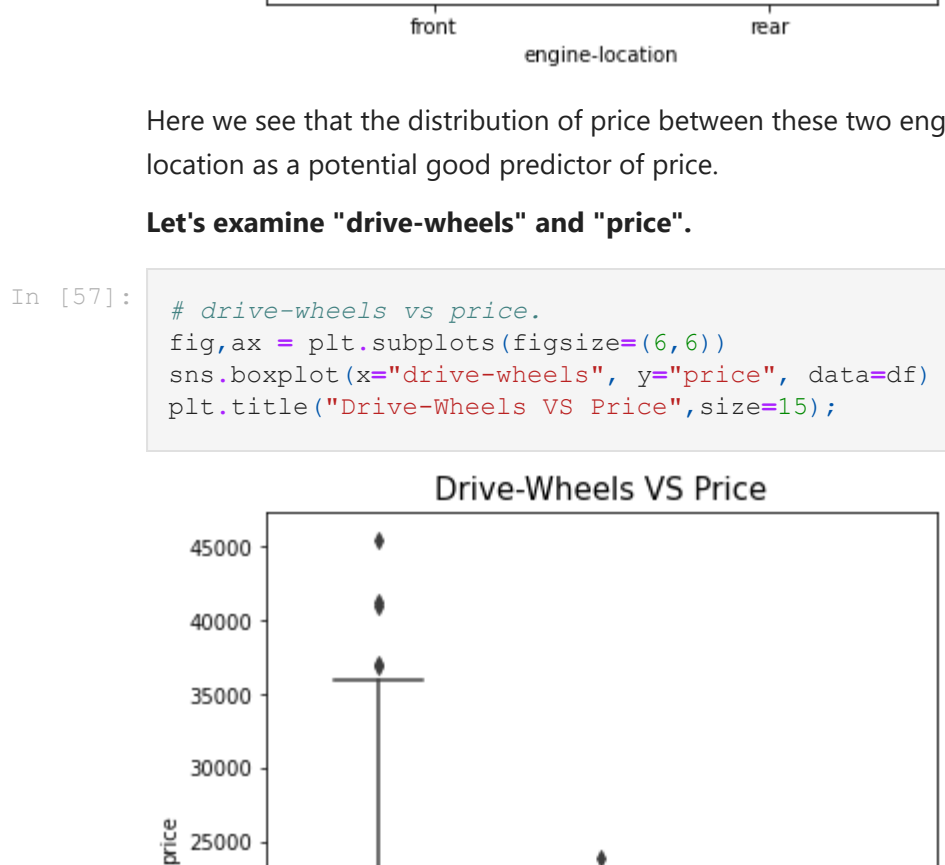
### Continuous Numerical Variables:

Continuous numerical variables are variables that may contain any value within some range. They can be of type 'int64' or 'float64'. A great way to visualize these variables is by using scatterplots with fitted lines.

In order to start understanding the **(linear) relationship** between an individual variable and the price, we can use "regplot" which plots the scatterplot plus the fitted regression line for the data.

#### Positive Linear Relationship

Let's find the **scatterplot** of "engine-size" and "price".



As the engine-size goes up, the price goes down; this indicates a positive direct correlation between these two variables.

Engine size seems like a pretty good predictor of price since the regression line is almost a perfect diagonal line.

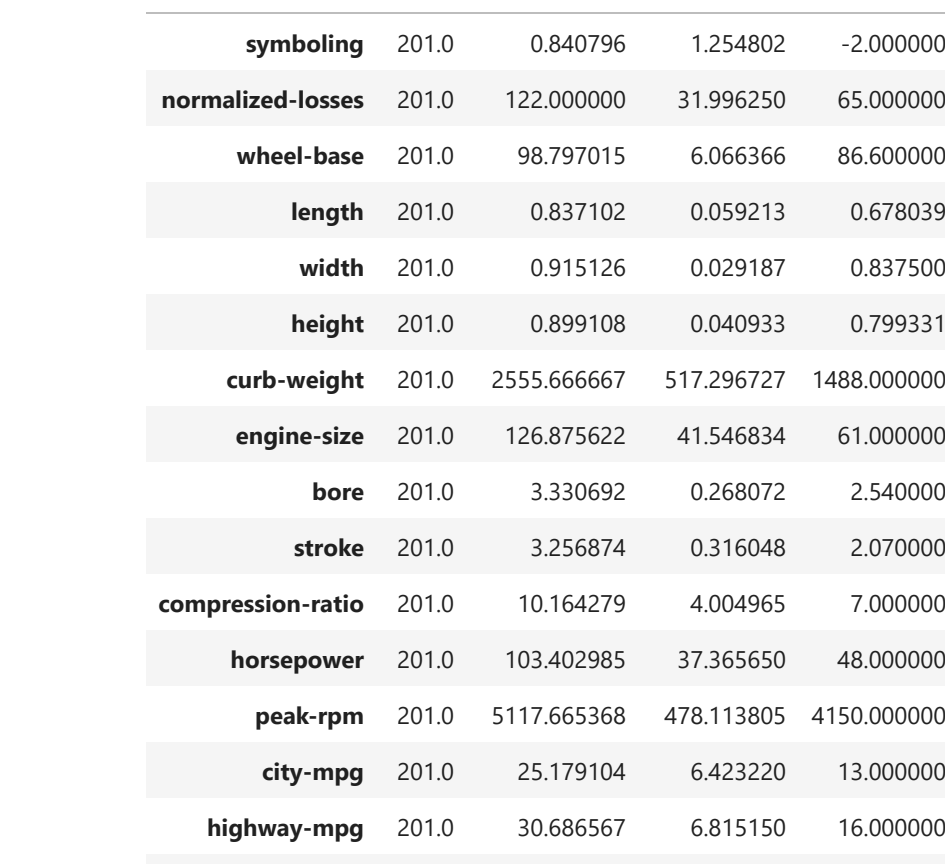
We can examine the correlation between 'engine-size' and 'price' and see that it's approximately 0.87.

```
In [48]: df[["engine-size", "price"]].corr()
```

	engine-size	price
engine-size	1.000000	0.872335
price	0.872335	1.000000

Highway-mpg is also a potential predictor variable of price.

**Let's find the scatterplot of "highway-mpg" and "price".**



As highway-mpg goes up, the price goes down; this indicates an inverse/negative relationship between these two variables.

Highway-mpg could potentially be a predictor of price.

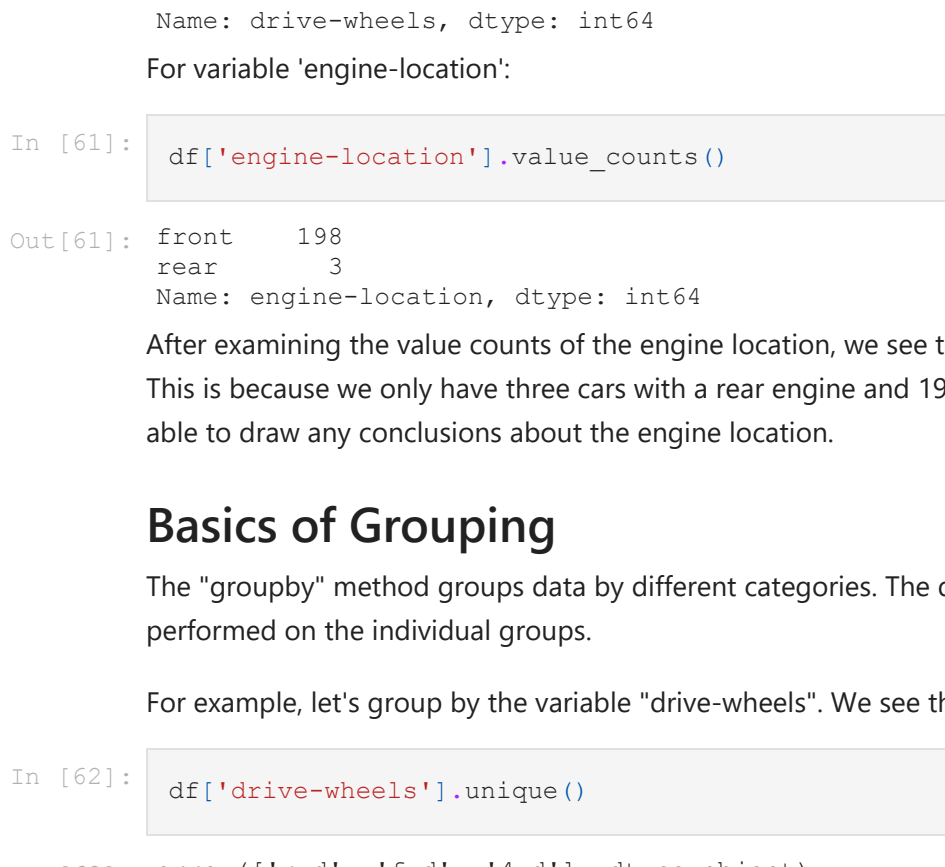
We can examine the correlation between 'highway-mpg' and 'price' and see it's approximately -0.704.

```
In [50]: df[["highway-mpg", "price"]].corr()
```

	highway-mpg	price
highway-mpg	1.000000	-0.704692
price	-0.704692	1.000000

#### Weak Linear Relationship

**Let's see if "peak-rpm" is a predictor variable of "price".**



Peak-rpm does not seem like a good predictor of the price at all since the regression line is close to horizontal.

Also, the data points are very scattered and far from the fitted line, showing lots of variability.

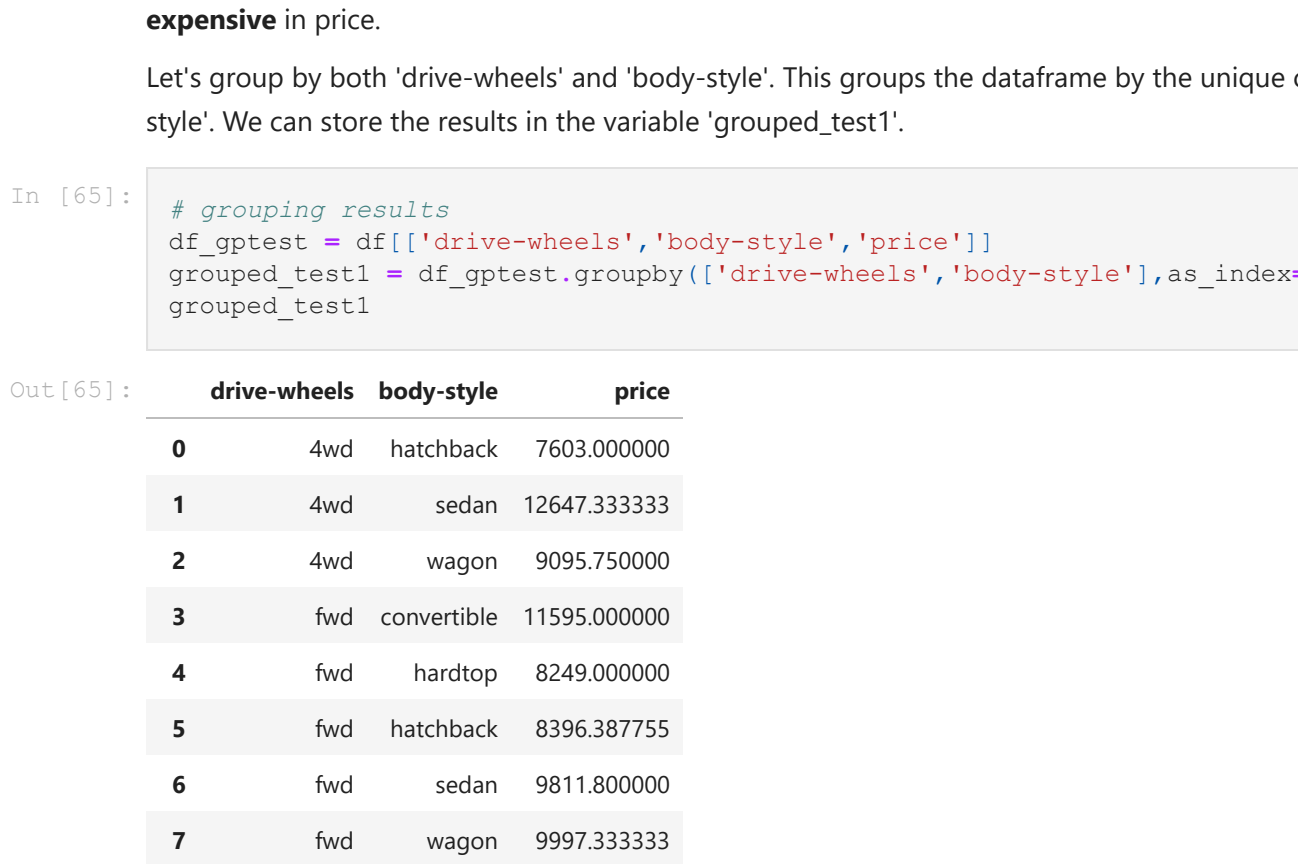
Therefore, it's not a reliable variable.

We can examine the correlation between 'peak-rpm' and 'price' and see it's approximately -0.101616.

```
In [52]: df[["peak-rpm", "price"]].corr()
```

	peak-rpm	price
peak-rpm	1.000000	-0.101616
price	-0.101616	1.000000

**Let's see if "stroke" is a predictor variable of "price".**



The correlation is 0.0823.

```
In [54]: df[["stroke", "price"]].corr()
```

	stroke	price
stroke	1.000000	0.082269
price	0.082269	1.000000

### Categorical Variables

These are variables that describe a 'characteristic' of a data unit.

A good way to visualize categorical variables is by using **boxplots**.

**Let's look at the relationship between "body-style" and "price".**



We see that the distributions of price between the different body-style categories have a significant overlap, so body-style cannot be a good predictor of price.

**Let's examine the relation between "engine-location" and "price".**



Here we see that the distribution of price between these two engine-location categories, front and rear, are distinct enough to take engine-location as a potential good predictor of price.

**Let's examine "drive-wheels" and "price".**



Here we see that the distribution of price between the different drive-wheels categories differs.

As such, drive-wheels could potentially be a predictor of price.

### Descriptive Statistical Analysis

Let's first take a look at the variables by utilizing a description method.

The `describe()` function automatically computes basic statistics for all continuous variables.

Any NaN values are automatically skipped in these statistics.

This will show:

```
In [58]: df.describe().T
```

	count	mean	std	min	25%	50%	75%	max
symboling	201.0	0.840796	1.254802	-2.000000	0.000000	1.000000	2.000000	3.00
normalized-losses	201.0	122.000000	31.995280	16.000000	65.000000	122.000000	137.000000	256.00
wheel-base	201.0	98.79						



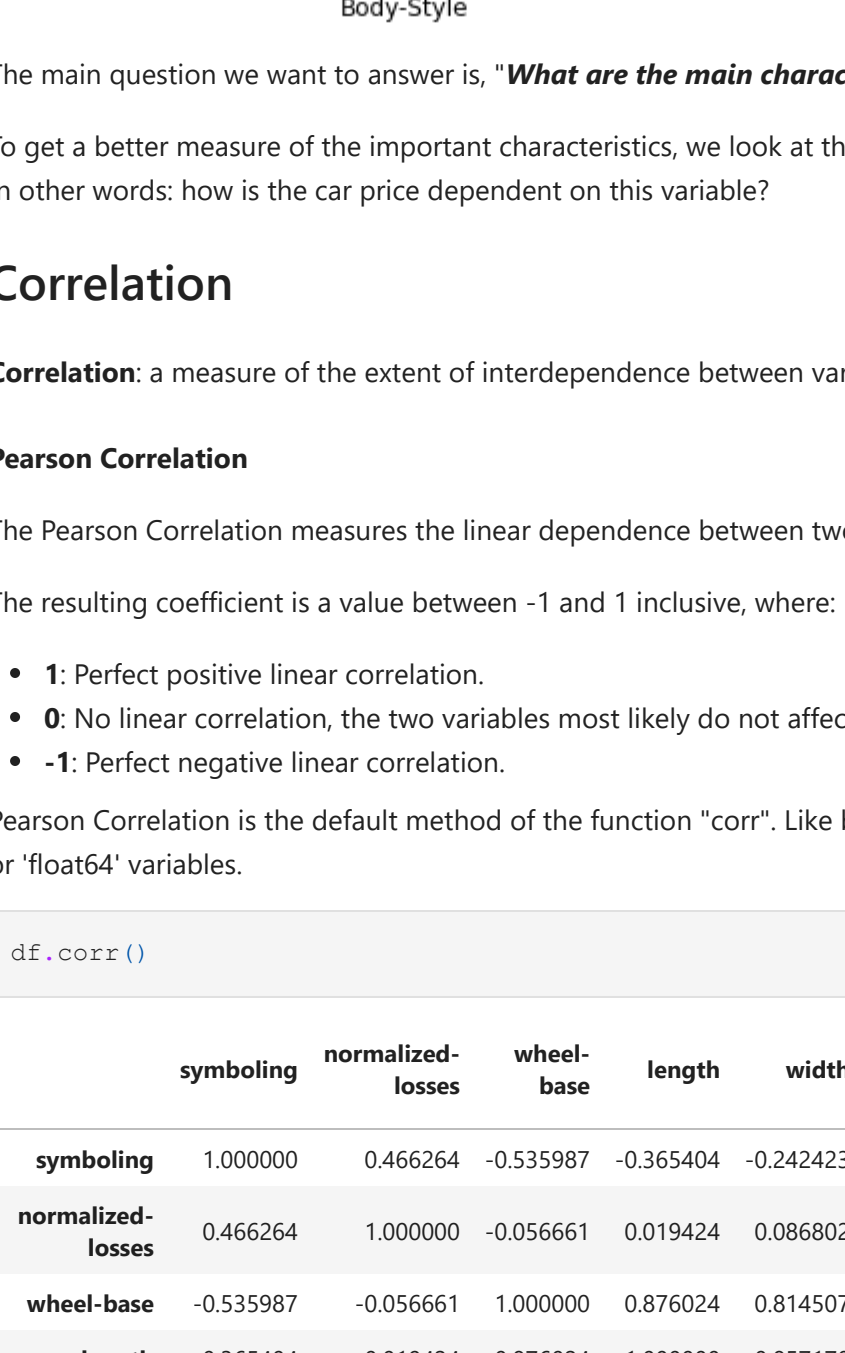
```
fig, ax = plt.subplots(figsize=(7,5))
im = ax.pcolor(grouped_pivot, cmap='b2b0u')

# Drive-Wheels & Labels to the center
ax.set_xticks(np.arange(grouped_pivot.shape[1]) + 0.5)
ax.set_yticks(np.arange(grouped_pivot.shape[0]) + 0.5)

# Insert Labels
ax.set_xticklabels(row_labels)
ax.set_yticklabels(col_labels)

# Rotate labels if too long
plt.xticks(rotation=45)

plt.title("Drive-Wheels & Body-Style VS Price",size=15)
plt.xlabel("Body-Style",size=12)
plt.ylabel("Drive-Wheels",size=12)
fig.colorbar(Im)
plt.show()
```



The main question we want to answer is, **"What are the main characteristics which have the most impact on the car price?"**.

To get a better measure of the important characteristics, we look at the correlation of these variables with the car price. In other words: how is the car price dependent on this variable?

## Correlation

**Correlation:** a measure of the extent of interdependence between variables. </b>

### Pearson Correlation

The Pearson Correlation measures the linear dependence between two variables X and Y.

The resulting coefficient is a value between -1 and 1 inclusive, where:

- 1: Perfect positive linear correlation.
- 0: No linear correlation, the two variables most likely do not affect each other.
- -1: Perfect negative linear correlation.

Pearson Correlation is the default method of the function "corr". Like before, we can calculate the Pearson Correlation of the of the 'mt54' or 'float64' variables.

```
In [68]: df.corr()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower	peak-rpm	city-mpg	highway-mpg	price
symboling	1.000000	0.466264	-0.535987	-0.365404	-0.242423	-0.550160	-0.233118	-0.110581	-0.140019	-0.008153	-0.182196	-0.081530	-0.279740	-0.352527	-0.424243	0.466264
normalized-losses	0.466264	1.000000	-0.056661	0.019424	0.086802	-0.373737	0.099404	0.112360	-0.029862	0.055045	-0.114713	-0.081530	-0.279740	-0.352527	-0.424243	0.466264
wheel-base	-0.535987	-0.056661	1.000000	0.076024	0.014507	0.590742	0.782097	0.572027	0.493244	0.158018	0.250313	0.158018	0.250313	0.158018	0.250313	0.466264
length	-0.365404	0.019424	0.076024	1.000000	0.085710	0.492063	0.880665	0.685025	0.608971	0.123952	0.123952	0.123952	0.123952	0.123952	0.123952	0.466264
width	-0.242423	0.086802	0.014507	0.085710	1.000000	0.306002	0.866201	0.729436	0.544885	0.188822	0.188822	0.188822	0.188822	0.188822	0.188822	0.466264
height	-0.550160	-0.373737	0.590742	0.492063	0.306002	1.000000	0.307581	0.074694	0.180449	-0.060663	0.259737	0.259737	0.259737	0.259737	0.259737	0.466264
curb-weight	-0.233118	0.099404	0.782097	0.880665	0.866201	0.307581	1.000000	0.849072	0.644060	0.167438	0.156433	0.156433	0.156433	0.156433	0.156433	0.466264
engine-size	-0.110581	0.112360	0.572027	0.685025	0.729436	0.074694	0.849072	1.000000	0.572609	0.205928	0.023889	0.023889	0.023889	0.023889	0.023889	0.466264
bore	-0.140019	-0.029862	0.493244	0.608971	0.544885	0.180449	0.644060	0.572609	1.000000	-0.055390	0.001263	0.001263	0.001263	0.001263	0.001263	0.466264
stroke	-0.008153	0.055045	0.158018	0.123952	0.188822	-0.060663	0.167438	0.205928	-0.055390	1.000000	0.187871	0.187871	0.187871	0.187871	0.187871	0.466264
compression-ratio	-0.182196	-0.114713	0.250313	0.159733	0.189867	0.259737	0.156433	0.023889	0.001263	0.187871	1.000000	0.187871	0.187871	0.187871	0.187871	0.466264
horsepower	-0.075810	0.217300	0.371718	0.579795	0.615056	-0.087001	0.757981	0.822668	0.566903	0.098128	-0.021489	0.021489	0.021489	0.021489	0.021489	0.466264
peak-rpm	-0.279740	-0.279740	-0.360305	-0.285970	-0.245800	-0.309974	-0.279361	-0.256733	-0.267392	-0.063561	-0.435780	-0.435780	-0.435780	-0.435780	-0.435780	0.466264
city-mpg	-0.352527	-0.352527	-0.470606	-0.470606	-0.633531	-0.048800	-0.749543	-0.650546	-0.582027	-0.033956	0.331425	0.331425	0.331425	0.331425	0.331425	0.466264
highway-mpg	-0.424243	-0.424243	-0.181877	-0.543034	-0.698142	-0.680635	-0.104812	-0.794889	-0.679571	-0.591309	-0.034636	0.268465	0.268465	0.268465	0.268465	0.466264
price	0.466264	0.466264	0.133999	0.584642	0.690628	0.751265	0.135486	0.834415	0.872335	0.543155	0.082269	0.071107	0.071107	0.071107	0.071107	0.466264
fuel-type-diesel	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	-0.082321	0.466264
fuel-type-gas	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.196735	0.466264

Sometimes we would like to know the significant of the correlation estimate.

### P-value

The P-value is the probability value that the correlation between these two variables is statistically significant.

Normally, we choose significance level of 0.05, which means we're 95% confident that correlation between the variables is significant.

By convention, when the:

- p-value is < 0.001: we say there is strong evidence that the correlation is significant.
- the p-value is < 0.05: there is moderate evidence that the correlation is significant.
- the p-value is < 0.1: there is weak evidence that the correlation is significant.
- the p-value is > 0.1: there is no evidence that the correlation is significant.

Also, Value of **e** is 2.71828182846.

We can obtain this information using 'stats' module in the 'scipy' library.

```
In [69]: from scipy import stats
```

## Wheel-Base vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'wheel-base' and 'price'.

```
In [70]: pearson_coef, p_value = stats.pearsonr(df['wheel-base'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.584641822655981 with a P-value of P = 8.076488270732989e-20
```

### Conclusion:

Since the p-value is < 0.001, the correlation between wheel-base and price is statistically significant, although the linear relationship isn't extremely strong (~0.585).

## Horsepower vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'horsepower' and 'price'.

```
In [71]: pearson_coef, p_value = stats.pearsonr(df['horsepower'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.8946068016571054 with a P-value of P = 6.273536270650504e-48
```

### Conclusion:

Since the p-value is < 0.001, the correlation between horsepower and price is statistically significant, and the linear relationship is quite strong (~0.899, close to 1).

## Length vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'length' and 'price'.

```
In [72]: pearson_coef, p_value = stats.pearsonr(df['length'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.6906283804483642 with a P-value of P = 8.016477466193759e-30
```

### Conclusion:

Since the p-value is < 0.001, the correlation between length and price is statistically significant, and the linear relationship is moderately strong (~0.691).

## Width vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'width' and 'price'.

```
In [73]: pearson_coef, p_value = stats.pearsonr(df['width'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.751265344522673 with a P-value of P = 9.200335510481646e-38
```

### Conclusion:

Since the p-value is < 0.001, the correlation between width and price is statistically significant, and the linear relationship is quite strong (~0.751).

## Curb-Weight vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'curb-weight' and 'price'.

```
In [74]: pearson_coef, p_value = stats.pearsonr(df['curb-weight'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.834415257702846 with a P-value of P = 2.189577238893691e-53
```

### Conclusion:

Since the p-value is < 0.001, the correlation between curb-weight and price is statistically significant, and the linear relationship is quite strong (~0.834).

## Engine-Size vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'engine-size' and 'price'.

```
In [75]: pearson_coef, p_value = stats.pearsonr(df['engine-size'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.8723351674453185 with a P-value of P = 9.26549162219389e-44
```

### Conclusion:

Since the p-value is < 0.001, the correlation between engine-size and price is statistically significant, and the linear relationship is very strong (~0.872).

## Bore vs. Price

Let's calculate the Pearson Correlation Coefficient and P-value of 'bore' and 'price'.

```
In [76]: pearson_coef, p_value = stats.pearsonr(df['bore'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is 0.5431553832626602 with a P-value of P = 8.04918948393489e-17
```

### Conclusion:

Since the p-value is < 0.001, the correlation between bore and price is statistically significant, but the linear relationship is only moderate (~0.521).

We can relate the process for each 'city-mpg' and 'highway-mpg':

## City-mpg vs. Price

```
In [77]: pearson_coef, p_value = stats.pearsonr(df['city-mpg'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is -0.68571067844677 with a P-value of P = 2.32113206567674e-29
```

### Conclusion:

Since the p-value is < 0.001, the correlation between city-mpg and price is statistically significant, and the coefficient of about -0.687 shows that the relationship is negative and moderately strong.

## Highway-mpg vs. Price

```
In [78]: pearson_coef, p_value = stats.pearsonr(df['highway-mpg'], df['price'])
print("The Pearson Correlation Coefficient is", pearson_coef, " with a P-value of P = ", p_value)

The Pearson Correlation Coefficient is -0.704692265089529 with a P-value of P = 1.749547114477352e-31
```

### Conclusion:

Since the p-value is < 0.001, the correlation between highway-mpg and price is statistically significant, and the coefficient of about -0.705 shows that the relationship is negative and moderately strong.

## ANOVA

Since ANOVA analyzes the difference between different groups of the same variable, the groupby function will come in handy. Because the ANOVA algorithm averages the data automatically, we do not need to take the average before hand.

## ANOVA: Analysis of Variance

The (ANOVA) is a statistical method used to test whether there are significant differences between the means of two or more groups.

ANOVA returns two parameters:

**F-test score:** It calculates the ratio of difference between the group's means over the variation within each of the sample group. A larger score means there is a larger difference between the means.

**P-value:** P-value tells how statistically significant our calculated score value is.

If our price variable is strongly correlated with the variable we are analyzing, we expect ANOVA to show a sizeable F-test score and a small p-value.

## Drive Wheels

```
In [79]: df.groupby('drive-wheels').price
```

drive-wheels	price
0	rwd 13495.0
1	rwd 16500.0
2	rwd hatchback 16500.0
3	fwd sedan 13950.0
4	4wd sedan 17450.0

201 rows x 3 columns

To see if different types of 'drive-wheels' impact 'price', we group the data.

```
In [80]: grouped_test2=df.groupby('drive-wheels').price
grouped_test2.head(2)
```

drive-wheels	price
0	rwd 13495.0
1	rwd 16500.0
2	rwd hatchback 16500.0
3	fwd sedan 13950.0
4	4wd sedan 17450.0

We can obtain the values of the method group using the method "get\_group".

```
In [81]: grouped_test2.get_group('4wd')['price']
```

4	17450.0
136	7603.0
140	9233.0
141	11259.0
144	8013.0
145	11694.0
150	7898.0
151	8778.0

We can use the function "f\_oneway" in the module 'stats' to obtain the **F-test score** and **P-value**.

```
In [82]: # ANOVA
f_val, p_val = stats.f_oneway(grouped_test2.get_group('fwd')['price'], grouped_test2.get_group('rwd')['price'],
                              grouped_test2.get_group('4wd')['price'])
print("ANOVA results: F=", f_val, ", P =", p_val)
```

ANOVA results: F= 67.95406500783399 , P = 3.3945443577151245e-23

This is a great result with a **large F-test score** showing a strong correlation and a P-value < 0.001 implying almost certain statistical significance. But does this mean all three tested groups are all this highly correlated?

Let's examine them separately.

### fwd and rwd

```
In [83]: f_val, p_val = stats.f_oneway(grouped_test2.get_group('fwd')['price'], grouped_test2.get_group('rwd')['price'])
print("ANOVA results: F=", f_val, ", P =", p_val)
```

ANOVA results: F= 130.553160959111 , P = 2.235536355677845e-23

**Conclusion:** They are highly correlated meaning, prices between "fwd" and "rwd" is significantly different as the F-score is large (130) and the p-value is < 0.001.

### 4wd and rwd

```
In [84]: f_val, p_val = stats.f_oneway(grouped_test2.get_group('4wd')['price'], grouped_test2.get_group('rwd')['price'])
print("ANOVA results: F=", f_val, ", P =", p_val)
```

ANOVA results: F= 8.580681368924756 , P = 0.00441149221225333

**Conclusion:** They are moderately correlated meaning, prices between "4wd" and "rwd" is moderately different as the F-score is 8 and the p-value is > 0.001 and is < 0.05.

### 4wd and fwd

```
In [85]: f_val, p_val = stats.f_oneway(grouped_test2.get_group('4wd')['price'], grouped_test2.get_group('fwd')['price'])
print("ANOVA results: F=", f_val, ", P =", p_val)
```

ANOVA results: F= 0.665465750252303 , P = 0.41620116697845666

### Conclusion:

They are weakly correlated meaning, prices between "4wd" and "fwd" are comparable as the F-score is small (0.6) and the p-value is > 0.05.

## Main Conclusion: Important Variables

We now have a better idea of our data and which variables are important to take into account when predicting the car price.

We have narrowed it down to the following variables:

### Continuous numerical variables:

- Length
- Width
- Curb-weight
- Engine-size
- Horsepower
- City-mpg
- Highway-mpg
- Wheel-base
- Bore

### Categorical variables:

- Drive-wheels

As we now move into building machine learning models to automate our analysis, feeding the model with variables that meaningfully affect our target variable will improve our model's prediction performance.

## Model Development

### Table of Contents

1. Linear Regression and Multiple Linear Regression
2. Model Evaluation Using Visualization
3. Polynomial Regression
4. Pipelines
5. Measures for Evaluation
6. Conclusion

Now, We will develop several models that will predict the price of the car using the variables or features. This is just an estimate but should give us an objective idea of how much the car should cost.

Some questions we want to ask now are:

- Do I know if the dealer is offering fair value for my trade-in?
- Do I know if I put a fair value on my car?

In data analytics, we use **Model Development** to help us predict future observations from the data we have.

A model will help us understand the exact relationship between different variables and how these variables are used to predict the result.

## Linear Regression and Multiple Linear Regression

### Linear Regression

One example of a Data Model that we will be using is:

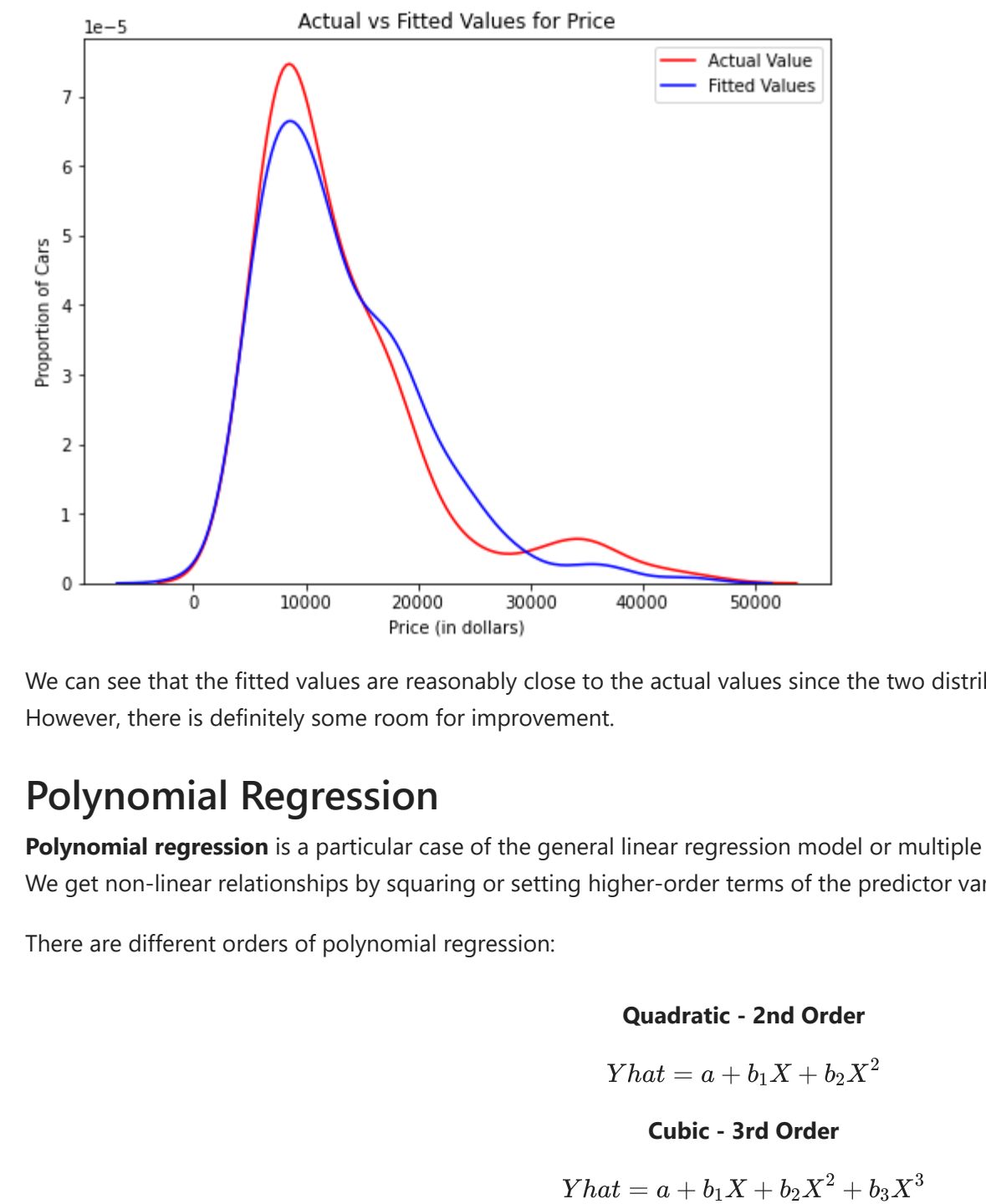
#### Simple Linear Regression

Simple Linear Regression is a method to help us understand the relationship between two variables:

- The predictor/independent variable (X)
- The response/dependent variable (that we want to predict)(Y)

The result of Linear Regression is a **linear function** that predicts the response (dependent)





We can see that the fitted values are reasonably close to the actual values since the two distributions overlap a bit. However, there is definitely some room for improvement.

## Polynomial Regression

**Polynomial regression** is a particular case of the general linear regression model or multiple linear regression models. We get non-linear relationships by squaring or setting higher-order terms of the predictor variables.

There are different orders of polynomial regression:

$$\text{Quadratic - 2nd Order} \\ \hat{y} = a + b_1X + b_2X^2$$

$$\text{Cubic - 3rd Order} \\ \hat{y} = a + b_1X + b_2X^2 + b_3X^3$$

$$\text{Higher-Order} \\ \hat{y} = a + b_1X + b_2X^2 + b_3X^3 + \dots$$

We saw earlier that a linear model did not provide the best fit while using "highway-mpg" as the predictor variable. Let's see if we can try fitting a polynomial model to the data instead.

We will use the following function to plot the data:

```
In [108]: def PlotPolly(model, independent_variable, dependent_variable, Name):
x_new = np.linspace(15, 55, 100)
y_new = model(x_new)

plt.figure(figsize=(8,6))
plt.plot(independent_variable, dependent_variable, '-', x_new, y_new, '-')
plt.title('Polynomial Fit for Price - %Name, size=15')
ax = plt.gca()
ax.set_facecolor((0.898, 0.898, 0.898))
fig = plt.gcf()
plt.ylim(0,)
plt.xlabel(Name)
plt.ylabel('Price of Cars')

plt.show()
plt.close()
```

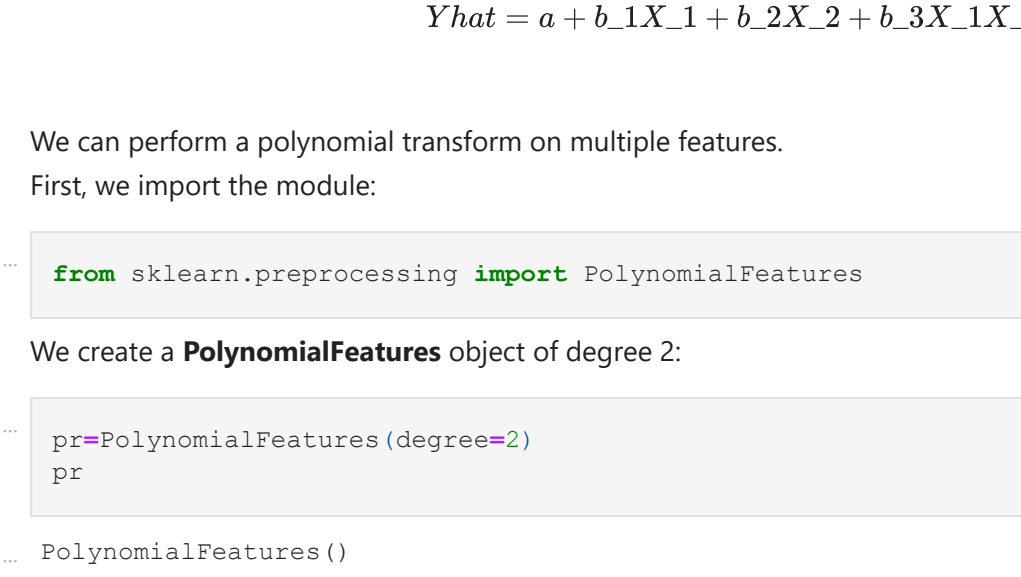
```
In [109]: # Let's get the variables:
x = df['highway-mpg']
y = df['price']
```

Let's fit the polynomial using the function **polyfit**, then use the function **poly1d** to display the polynomial function.

```
In [110]: # Here we use a polynomial of the 3rd order (cubic).
f = np.polyfit(x, y, 3)
p = np.poly1d(f)
print(p)

-1.557 x3 + 204.8 x2 - 8965 x + 1.379e+05
```

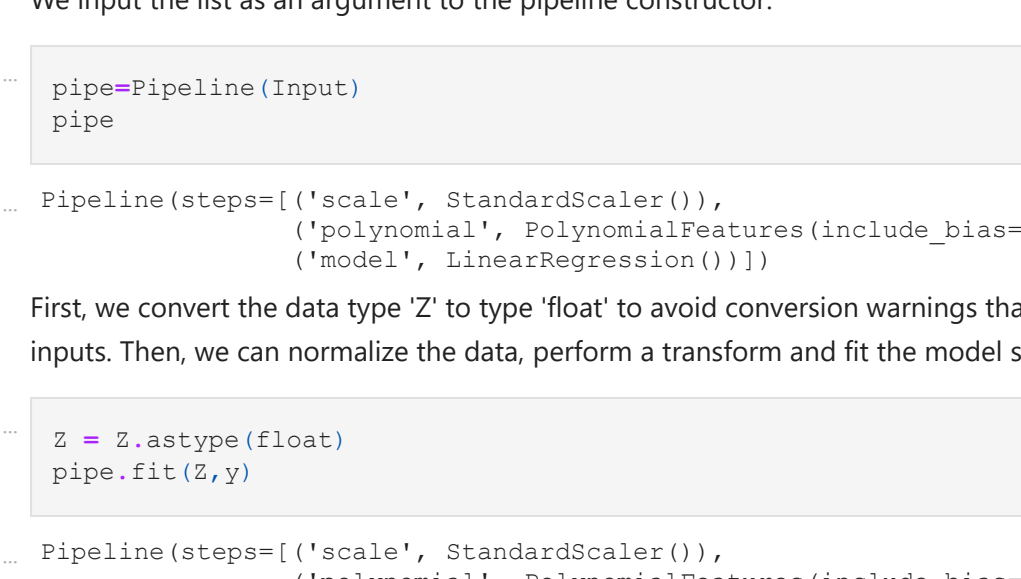
```
In [111]: # Let's plot the function:
PlotPolly(p, x, y, 'highway-mpg')
```



We can already see from plotting that this polynomial model performs better than the linear model. This is because the generated polynomial function "fits" more of the data points.

Let's try Creating 11 order polynomial model with the variables x and y from above.

```
In [112]: # Here we use a polynomial of the 11th order.
f1 = np.polyfit(x, y, 11)
p1 = np.poly1d(f1)
print(p1, "\n")
PlotPolly(p1, x, y, 'Highway-MPG')
```



The analytical expression for Multivariate Polynomial function gets complicated. For example, the expression for a second-order (degree=2) polynomial with two variables is given by:

$$\hat{y} = a + b_1X_1 + b_2X_2 + b_3X_1X_2 + b_4X_1^2 + b_5X_2^2$$

We can perform a polynomial transform on multiple features.

First, we import the module:

```
In [113]: from sklearn.preprocessing import PolynomialFeatures
```

We create a **PolynomialFeatures** object of degree 2:

```
In [114]: pr=PolynomialFeatures(degree=2)
pr
```

```
Out [114]: PolynomialFeatures()
```

```
In [115]: Z_pr=pr.fit_transform(Z)
```

In the original data, there are 201 samples and 4 features.

```
In [116]: Z.shape
```

```
Out [116]: (201, 4)
```

After the transformation, there are 201 samples and 15 features.

```
In [117]: Z_pr.shape
```

```
Out [117]: (201, 15)
```

## Pipeline

Data Pipelines simplify the steps of processing the data. We use the module **Pipeline** to create a pipeline. We also use **StandardScaler** as a part in our pipeline.

```
In [118]: from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
```

We create the pipeline by creating a list of tuples including the name of the model or estimator and its corresponding constructor.

```
In [119]: Input=[('scale',StandardScaler()), ('polynomial', PolynomialFeatures(degree=2,include_bias=False)), ('model',LinearRegression())]
```

We input the list as an argument to the pipeline constructor:

```
In [120]: pipe=Pipeline(Input)
pipe
```

```
Out [120]: Pipeline(steps=[('scale', StandardScaler()), ('polynomial', PolynomialFeatures(include_bias=False)), ('model', LinearRegression())])
```

First, we convert the data type 'Z' to type 'float' to avoid conversion warnings that may appear as a result of StandardScaler taking float inputs. Then, we can normalize the data, perform a transform and fit the model simultaneously.

```
In [121]: Z = Z.astype(float)
pipe.fit(Z,y)
```

```
Out [121]: Pipeline(steps=[('scale', StandardScaler()), ('polynomial', PolynomialFeatures(include_bias=False)), ('model', LinearRegression())])
```

Similarly, we can normalize the data, perform a transform and produce a prediction simultaneously.

```
In [122]: ypipe=pipe.predict(Z)
ypipe[0:4]
```

```
Out [122]: array([13102.93329646, 13102.93329646, 18226.43450275, 10391.09183955])
```

## Measures for Evaluation

When evaluating our models, not only do we want to visualize the results, but we also want a quantitative measure to determine how accurate the model is.

Two very important measures that are often used in Statistics to determine the accuracy of a model are:

- **R<sup>2</sup> / R-squared**
- **Mean Squared Error (MSE)**

### R-squared

**R-squared**, also known as the **coefficient of determination**, is a measure to indicate how close the data is to the fitted regression line. The value of the R-squared is the percentage of variation of the response variable (y) that is explained by a linear model.

### Mean Squared Error (MSE)

The **Mean Squared Error** measures the average of the squares of errors. That is, the difference between actual value (y) and the estimated value (ŷ).

```
In [123]: # Creating a Dictionary for Storing R-Squared Values of Different Models.
evalDict={}

# Let's calculate the R^2:
```

```
In [124]: # highway_mpg fit
lm_fit(x, y)
# Find the R^2
r2 = lm.score(x, y)
print("The R-square is: ", r2)
evalDict['SLR'] = r2 * 100

The R-square is: 0.4965918843391759
```

So we can say that **~49.659%** of the variation of the price is explained by this simple linear model "highway\_mpg\_fit".

Let's calculate the MSE:

```
In [125]: We can predict the output i.e., "yhat" using the predict(), where X is the input variable:

Yhat=lm.predict(X)
print("The output of the first four predicted value is: ", Yhat[0:4])

The output of the first four predicted value is: [16236.50464347 16236.50464347 17058.23802179 13771.3045085 ]
```

```
In [126]: # Import the function "mean_squared_error" from the module "metrics":
from sklearn.metrics import mean_squared_error
```

We can compare the predicted results with the actual results:

```
In [127]: mse = mean_squared_error(df['price'], Yhat)
print("The mean square error of price and predicted value is: ", mse)

The mean square error of price and predicted value is: 31635042.94639895
```

## Model 2: Multiple Linear Regression

Let's calculate the R<sup>2</sup>:

```
In [128]: # fit the model
lm_fit(Z, df['price'])
# Find the R^2
r2 = lm.score(Z, df['price'])
print("The R-square is: ", r2)
evalDict['MLR'] = r2 * 100

The R-square is: 0.8093732522175299
```

We can say that **~80.937%** of the variation of the price is explained by this multiple linear regression "multi\_fit".

Let's calculate the MSE:

We produce a prediction:

```
In [129]: y_predict_multifit = lm.predict(Z)
```

We compare the predicted results with the actual results:

```
In [130]: print("The mean square error of price and predicted value using multifit is:",
mean_squared_error(df['price'], y_predict_multifit))

The mean square error of price and predicted value using multifit is: 11979300.349818885
```

## Model 3: Polynomial Fit

Let's calculate the R<sup>2</sup>:

Let's import the function **r2\_score** from the module **metrics** as we are using a different function.

```
In [131]: from sklearn.metrics import r2_score

We apply the function to get the value of R^2:
```

```
In [132]: r2 = r2_score(y, p(x))
print("The R-square value is: ", r2)
evalDict['Polynomial Reg.'] = r2*100

The R-square value is: 0.6741946663906513
```

We can say that **~67.419%** of the variation of the price is explained by this polynomial fit.

We can also calculate the MSE:

```
In [133]: print("The mean square error of price and predicted value is:",
mean_squared_error(df['price'], p(x)) )

The mean square error of price and predicted value is: 20474146.42636125
```

## Prediction and Decision Making

### Decision Making: Determining a Good Model Fit

Now that we have Visualized the different models, and generated the R-squared and MSE values for the fits, how do we determine a good model fit?

- **What is a good R-squared value?**

When comparing models, **the model with the higher R-squared value is a better fit** for the data.

- **What is a good MSE?**

When comparing models, **the model with the smallest MSE value is a better fit** for the data.

Let's take a look at the values for these models.

**Simple Linear Regression (SLR):** Using Highway-mpg as a Predictor Variable of Price.

- R-squared: 0.4965918843391759
- MSE: 3.16 x10<sup>7</sup>

**Multiple Linear Regression (MLR):** Using Horsepower, Curb-weight, Engine-size, and Highway-mpg as Predictor Variables of Price.

- R-squared: 0.80896354913783497
- MSE: 1.2 x10<sup>7</sup>

**Polynomial Fit:** Using Highway-mpg as a Predictor Variable of Price.

- R-squared: 0.6741946663906514
- MSE: 2.05 x 10<sup>7</sup>

### (SLR) Model vs (MLR) Model

Sometimes, the more variables you have, the better your model is at predicting, but this is not always true.

Sometimes we may not have enough data, you may run into numerical problems, or many of the variables may not be useful and even act as noise.

In order to compare the results of the MLR vs SLR models, we look at a combination of both the R-squared and MSE to make the best conclusion about the model.

- **MSE:** The MSE of SLR is **3.16x10<sup>7</sup>** while MLR has an MSE of **1.2 x10<sup>7</sup>**. The MSE of MLR is much smaller.
- **R-squared:** In this case, we can also see that there is a big difference between the R-squared of the SLR and the R-squared of the MLR. The R-squared for the SLR (**-0.497**) is very small compared to the R-squared for the MLR (**-0.809**).

This R-squared in combination with the MSE show that MLR seems like the better model fit in this case compared to SLR.

### (SLR) vs. Polynomial Fit

- **MSE:** We can see that Polynomial Fit brought down the MSE, since its MSE (**2.05x10<sup>7</sup>**) is smaller than the one from the SLR (**3.16x10<sup>7</sup>**).
- **R-squared:** The R-squared for the Polynomial Fit (**-0.674**) is larger than the R-squared for the SLR (**-0.497**), so the Polynomial Fit also brought up the R-squared quite a bit.

Since the Polynomial Fit resulted in a lower MSE and a higher R-squared, we can conclude that this was a better fit model than the simple linear regression for predicting "price".

## Multiple Linear Regression (MLR) vs. Polynomial Fit

- **MSE:** The MSE for the MLR is smaller than the MSE for the Polynomial Fit.
- **R-squared:** The R-squared for the MLR is also much larger than for the Polynomial Fit.

## Conclusion

Comparing these three models, we conclude that **the MLR model is the best model** to be able to predict price from our dataset. This result makes sense since we have 26 variables in total and we know that more than one of those variables are potential predictors of the final car price.

## Model Evaluation and Refinement

We have built models and made predictions of vehicle prices.

Now we will determine how accurate these predictions are.

## Table of Contents

- [Model Evaluation](#)
- [Over-fitting and Model Selection](#)
- [Ridge Regression](#)
- [Grid Search](#)
- [Model Success Plot](#)

```
In [136]: # Library for interactive plotting.
from IPython.display import interact
```

First lets only use numeric data.

```
In [137]: df=df.get_numeric_data()
df.head()
```

	symboling	normalized-losses	wheel-base	length	width	height	curb-weight	engine-size	bore	stroke	compression-ratio	horsepower	peak-rpm	city-mpg	highway-mpg
0	3	122	88.6	0.811148	0.890278	0.816054	2548	130	3.47	2.68	9.0	111	5000.0	21	
1	3	122	88.6	0.811148	0.890278	0.816054	2548	130	3.47	2.68	9.0	111	5000.0	21	
2	1	122	94.5	0.822681	0.909722	0.876254	2823	152	2.68	3.47	9.0	154	5000.0	19	
3	2	164	99.8	0.848630	0.919444	0.908027	2337	109	3.19	3.40	10.0	102	5500.0	18	
4	2	164	99.4	0.848630	0.922222	0.908027	2824	136	3.19	3.40	8.0	115	5500.0	24	

## User-Defined Functions for plotting

```
In [138]: def DistributionPlot(x_train, BlueFunction, RedName, BlueName, Title=""):
width = 8
height = 6
plt.figure(figsize=(width, height))

ax1 = sns.kdeplot(x=x_train, color='b', label=RedName)
ax2 = sns.kdeplot(x=x_train, color='r', label=BlueName, ax=ax1)

plt.title(Title, size=15)
plt.xlabel('Price (in dollars)', size=13)
plt.ylabel('Proportion of Cars', size=13)
plt.legend()

plt.show()
plt.close()
```

```
In [139]: def PollyPlot(x_train, x_test, y_train, y_test, lr, poly_transform, text=""):
width = 8
height = 6
plt.figure(figsize=(width, height))

#training data
# lr: linear regression object
# poly_transform: polynomial transformation object

xmax=max(x_train.values.max(), x_test.values.max())
xmin=min(x_train.values.min(), x_test.values.min())
xmp=arange(xmin, xmax, 0.1))

plt.plot(x_train, y_train, 'co', label='Training Data')
plt.plot(x_test, y_test, 'go', label='Test Data')
plt.plot(x, lr.predict(poly_transform.fit_transform(x.reshape(-1, 1))), label='Predicted Function')
plt.ylim((0, 60000))
plt.ylabel('Price', size=15)
plt.xlabel('Price of cars', size=13)
plt.legend()

plt.show()
plt.close()
```

## Training and Testing

An important step in testing the 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 [139]: y_data = df['price']
```

Drop price data in **x\_data**

```
In [140]: 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 [140]: from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.15, random_state=0)

print("Number of test samples : ", x_test.shape[0])
print("Number of training samples: ", x_train.shape[0])

Number of test samples : 31
Number of training samples: 170

The test_size parameter sets the proportion of data that is split into the testing set.
In the above, the testing set is set to 15% of the total dataset.

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

```
In [141]: from sklearn.linear_model import LinearRegression
```

We create a Linear Regression object:

```
In [142]: lre=LinearRegression()
```

We fit the model using the feature horsepower

```
In [143]: lre.fit(x_train[['horsepower']], y_train)
```

```
Out [143]: LinearRegression()
```

Let's Calculate the R<sup>2</sup> on the test data:

```
In [144]: lre.score(x_test[['horsepower']], y_test)
```

```
Out [144]: 0.7076967079117261
```

We can see the R<sup>2</sup> is smaller using the train data.

```
In [145]: lre.score(x_train[['horsepower']], y_train)
```

```
Out [145]: 0.6450110239384648
```

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

## Cross-validation Score

Lets import **model\_selection** from the module **cross\_val\_score**.

```
In [146]: from sklearn.model_selection import cross_val_score
```

We input the object, the feature in this case **horsepower**, the target data (**y\_data**). The parameter **cv** determines the number of folds; in this case 4.

```
In [147]: kcross = 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 in the fold:
```

```
In [148]: kcross
array([0.7465419, 0.51718424, 0.74814454, 0.04825398])
```

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

```
In [149]: print("The mean of the folds is : ", round(kcross.mean(),3))
print("The standard deviation is : ", round(kcross.std(),3))

The mean of the folds is : 0.522
The standard deviation is : 0.291
```

We can use **mean\_squared\_error** as a score by setting the parameter **scoring** metric to **neg\_mean\_squared\_error**.

```
In [150]: -1 * cross_val_score(lre, x_data[['horsepower']], y_data, cv=4, scoring='neg_mean_squared_error')

array([120251357.7835463, 43743920.05390439, 12525158.34507633,
17564357.8926653])
```

We can also use the function **cross\_val\_predict** to predict the output.

The function splits up the data into the specified number of folds, using one fold to get a prediction while the rest of the folds are used as training data.

First import the function:

```
In [151]: from sklearn.model_selection import cross_val_predict
```

We input the model object, the feature in this case **'horsepower'**, the target data **y\_data**. Then, we can produce an output:

```
In [152]: yhat = cross_val_predict(lre, x_data[['horsepower']], y_data, cv=4)
yhat[0:5]
```

```
Out [152]: array([14142.23793549, 14142.23793549, 20815.3029844, 12745.549902,
14762.9881726 ])
```

## Overfitting 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.

Let's create Multiple linear regression objects and train the model using "horsepower", "curb-weight", "engine-size" and "highway-mpg" as features.

```
In [153]: lr = LinearRegression()
lr.fit(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']], y_train)
```

```
Out [153]: LinearRegression()
```

Prediction using training data:

```
In [154]: yhat_train = lr.predict(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
yhat_train[0:5]
```

```
Out [154]: array([11927.25153792, 11236.70125955, 6436.82274615, 21891.09877661,
16682.10139521])
```

Prediction using test data:

```
In [155]: yhat_test = lr.predict(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
yhat_test[0:5]
```

```
Out [155]: array([11349.48964574, 5914.6130239, 11243.35261505, 6661.95904136,
15555.34734434])
```

Let's perform some model **evaluation** using our training and testing data separately.

First we import the **seaborn** and **matplotlib** library for plotting.

```
In [156]: import matplotlib.pyplot as plt
import seaborn as sns
```

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

```
In [157]: Title = "Distribution Plot of Predicted vs Actual Values (Train Data)"
DistributionPlot(y_train, yhat_train, "Actual Values (Train)", "Predicted Values (Train)", Title)
```

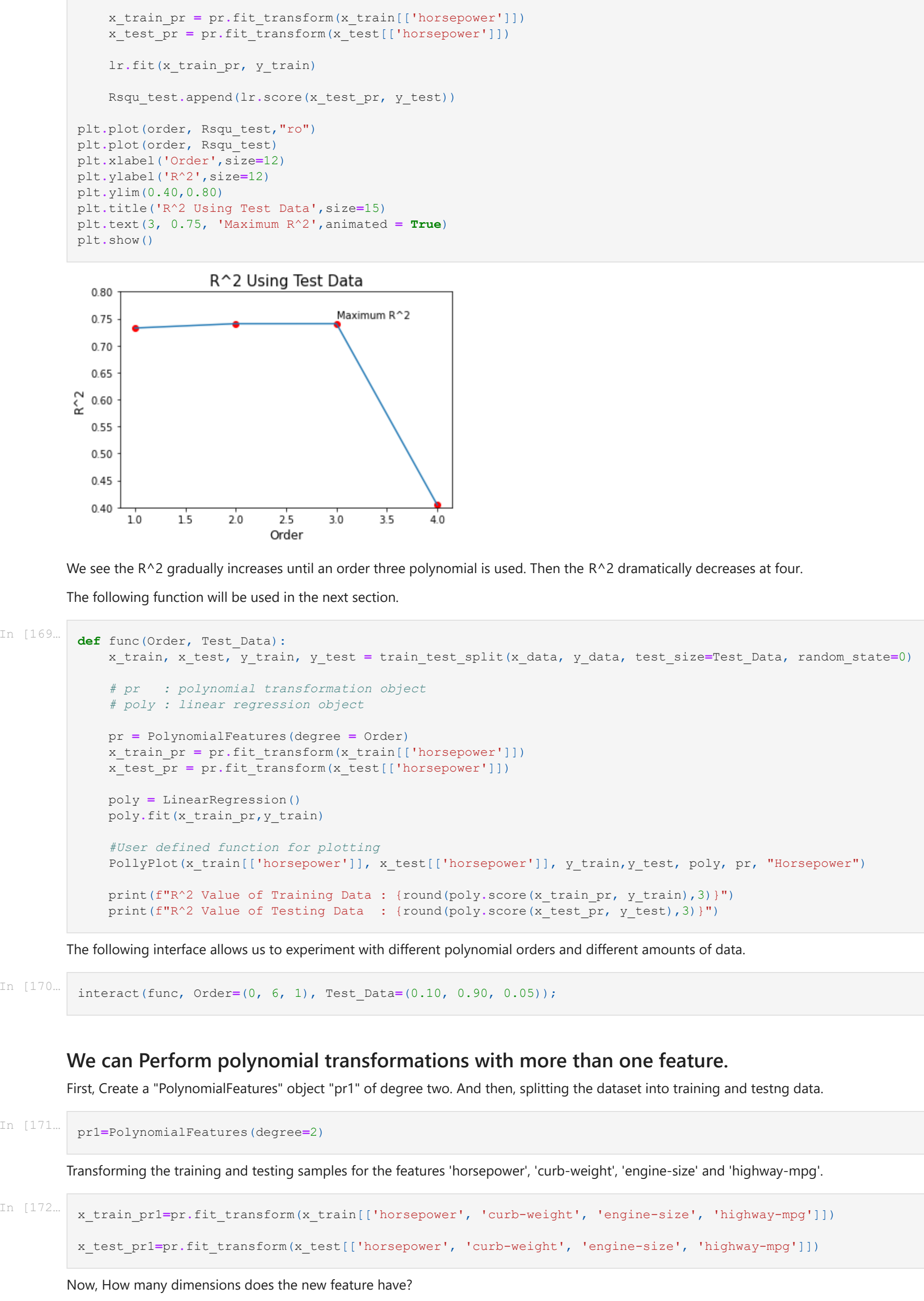




```
In [167]: poly.score(x_test_pr, y_test)
Out[167]: -29.815481910638443
```

We see the R<sup>2</sup> for the training data is 0.5567 while the R<sup>2</sup> on the test data was -29.87.  
The lower the R<sup>2</sup>, the worse the model; a Negative R<sup>2</sup> is a sign of Overfitting.

Let's see how the R<sup>2</sup> changes on test data for different Order Polynomials and Plot the results:



We see the R<sup>2</sup> gradually increases until an order three polynomial is used. Then the R<sup>2</sup> dramatically decreases after four.  
The following function will be used in the next section.

```
In [169]: def func(Order, TestData):
x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=TestData, random_state=0)
# pr : polynomial transformation object
# poly : linear regression object
pr = PolynomialFeatures(degree = Order)
x_train_pr = pr.fit_transform(x_train[['horsepower']])
x_test_pr = pr.fit_transform(x_test[['horsepower']])
poly = LinearRegression()
poly.fit(x_train_pr, y_train)
#User defined function for plotting
PolyPlot(x_train[['horsepower']], x_test[['horsepower']], y_train, y_test, poly, pr, "Horsepower")
print(f'R^2 Value of Training Data : {round(poly.score(x_train_pr, y_train),3)}')
print(f'R^2 Value of Testing Data : {round(poly.score(x_test_pr, y_test),3)}')
```

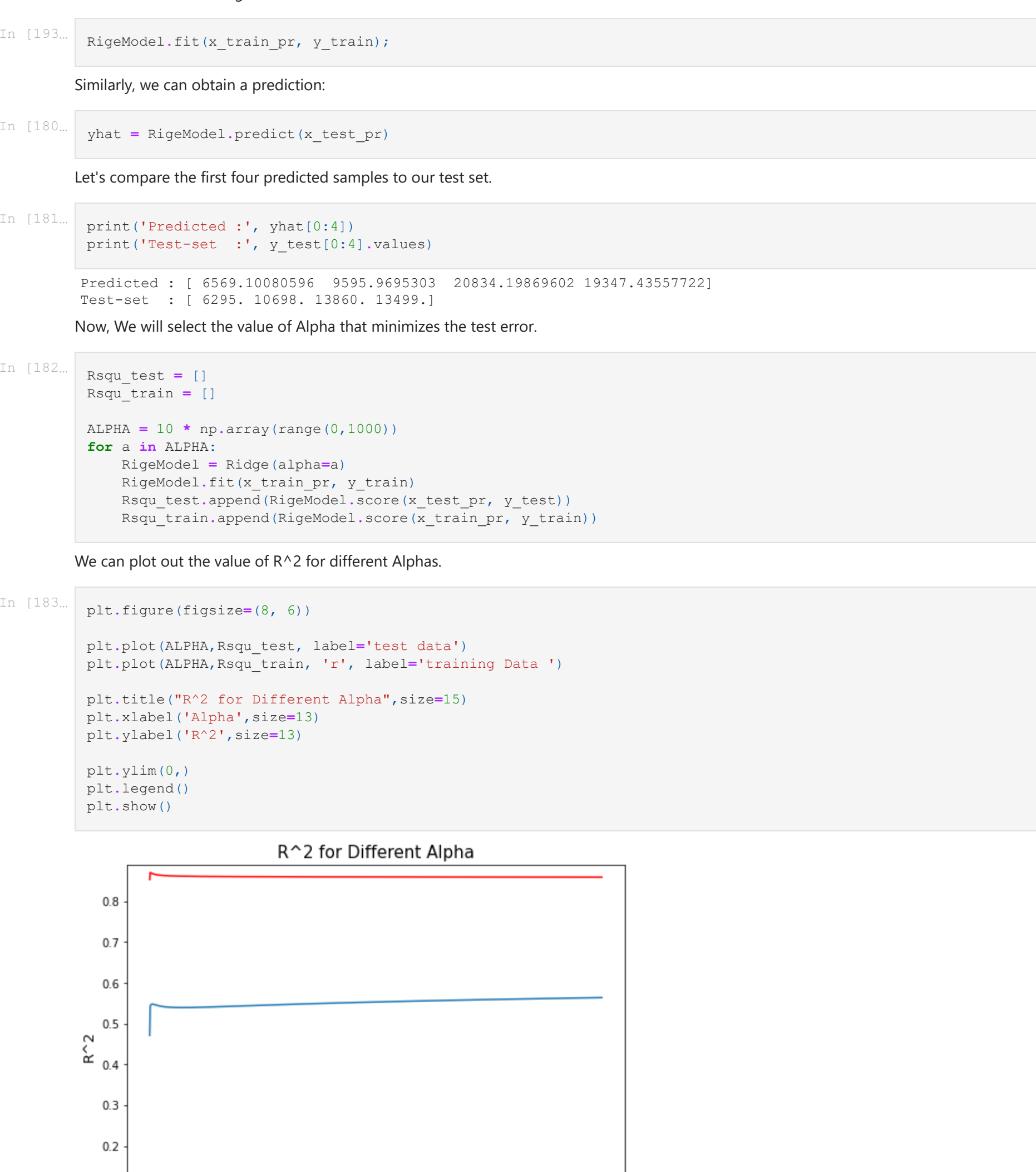
The following interface allows us to experiment with different polynomial orders and different amounts of data.

```
In [170]: interact(func, Order=(0, 6, 1), Test_Data=(0.10, 0.90, 0.05)):
```

We can Perform polynomial transformations with more than one feature.  
First, Create a "PolynomialFeatures" object "pr" of degree two. And then, splitting the dataset into training and testing data.

```
In [171]: pr=PolynomialFeatures(degree=2)
Transforming the training and testing samples for the features 'horsepower', 'curb-weight', 'engine-size' and 'highway-mpg'.
In [172]: x_train_pr=pr.fit_transform(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
x_test_pr=pr.fit_transform(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg']])
Now, How many dimensions does the new feature have?
In [173]: x_train_pr.shape
Out[173]: (110, 70)
Creating a linear regression model "poly1" and training the object.
In [174]: poly1 = LinearRegression().fit(x_train_pr, y_train)
```

Predicting an output on the polynomial features, then using the function "DistributionPlot" to display the distribution of the predicted output vs the test data.



**Conclusion:** The predicted value is higher than actual value for cars where the price \$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.

## Ridge regression

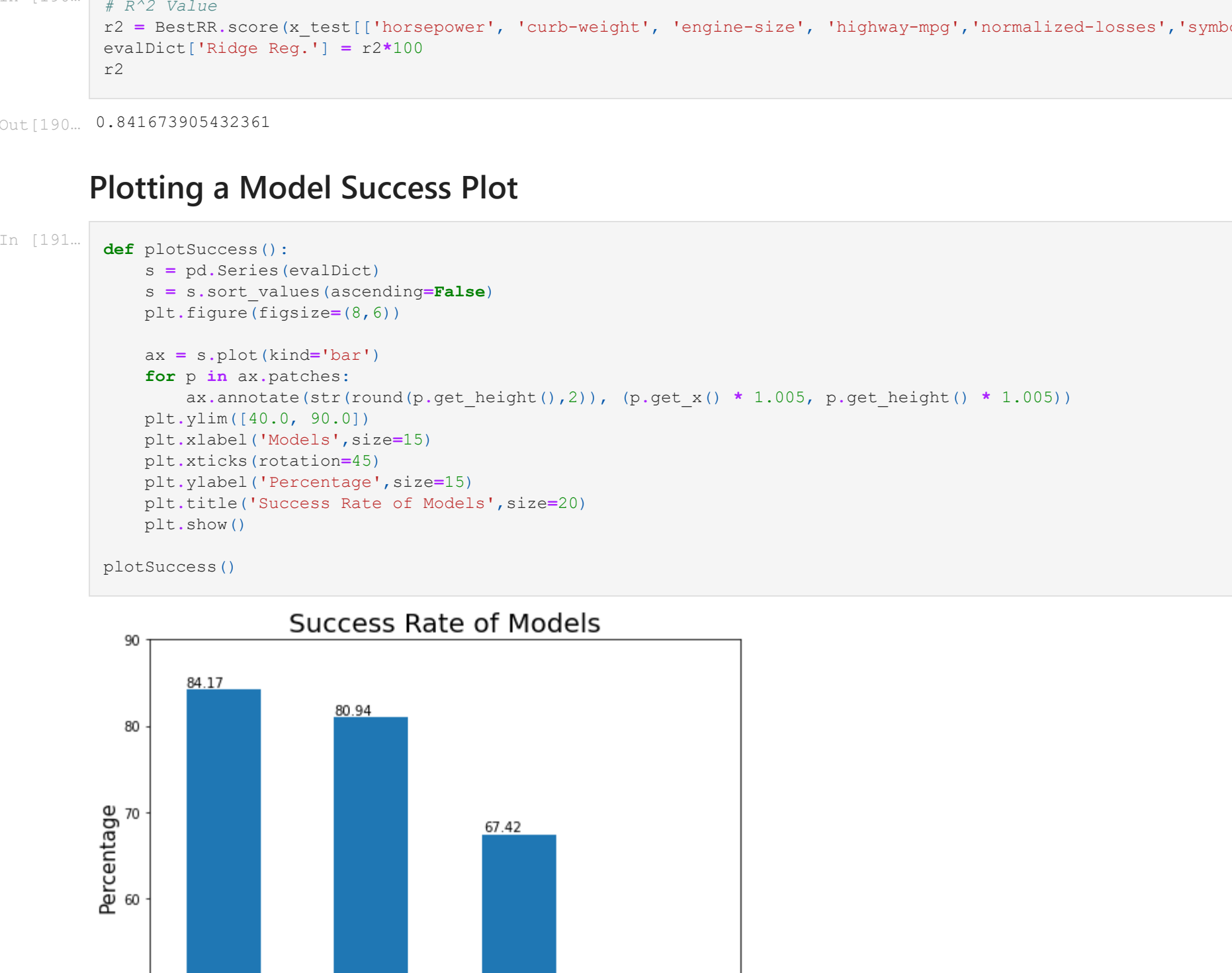
**Ridge Regression** is usually used when there is a high correlation between independent variables. This is because, in the case of multi collinear data, the least square estimates give unbiased values. This is a powerful regression method where the model is less susceptible to overfitting.

It performs a degree two polynomial transformation on our data.

```
In [176]: pr=PolynomialFeatures(degree=2)
x_train_pr=pr.fit_transform(x_train[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg'], 'normalized-losses'])
x_test_pr=pr.fit_transform(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg'], 'normalized-losses'])
Let's import Ridge from the module linear models.
In [177]: from sklearn.linear_model import Ridge
Let's create a Ridge regression object, setting the Regularization Parameter "alpha" to 0.1.
RidgeModel=Ridge(alpha=0.1)
Now, fit the model using the method fit.
In [178]: RidgeModel.fit(x_train_pr, y_train)
```

Similarly, we can obtain a prediction:

```
In [180]: yhat = RidgeModel.predict(x_test_pr)
Let's compare the first four predicted samples to our test set.
In [181]: print('Predicted :', yhat[0:4])
print('Test-set :', y_test[0:4].values)
Predicted : [ 6569.10080596  9595.9695303  20834.19869602 19347.43557722]
Test-set : [ 6295. 10698. 13860. 13499.]
Now, We will select the value of Alpha that minimizes the test error.
```



The red line in figure represents the R<sup>2</sup> of the training data, as Alpha increases the R<sup>2</sup> decreases; therefore as Alpha increases the model performs worse on the training data.  
The blue line represents the R<sup>2</sup> on the test data, as the value for Alpha increases the R<sup>2</sup> increases.

## 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 [184]: from sklearn.model_selection import GridSearchCV
We create a dictionary of parameter values:
In [185]: parameters1 = {'alpha': [0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000, 100000], 'normalize': [True, False]}
Creating a ridge regression object:
In [186]: RR=Ridge()
Out[186]: Ridge()
```

Creating a Ridge grid-search object.

```
In [187]: Grid1 = GridSearchCV(RR, parameters1, cv=4)
Fit the model
In [188]: Grid1.fit(x_data[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg'], 'normalized-losses'], 'symboling')],
Out[188]: GridSearchCV(cv=4, estimator=Ridge(),
param_grid=[('alpha': [0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000,
100000],
'normalize': [True, False])])
```

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 [189]: BestRR=Grid1.best_estimator_
BestRR
Out[189]: Ridge(alpha=0.1, normalize=True)
We now test our model on the test data
In [190]: # R^2 Value
r2 = BestRR.score(x_test[['horsepower', 'curb-weight', 'engine-size', 'highway-mpg'], 'normalized-losses'], 'symboling')
evalDict['Ridge Reg.' ] = r2*100
Out[190]: 0.841673905432361
```

## Plotting a Model Success Plot



## Conclusion:

Ridge Regression Model has the Highest Success Rate. So, we can use it to predict the prices of Cars successfully.  
A Test Case is presented below.

```
In [192]: # Real Life Scenario
# Variables can have any appropriate Value.
# Let's say your friend's car has the following features.
horsepower = 110.0
curb_weight = 2350.0
engine_size = 120.0
highway_mpg = 25.0
normalized_losses = 110.0
symboling = 1.0
testCase = [[horsepower, curb_weight, engine_size, highway_mpg, normalized_losses, symboling]]
testDf = pd.DataFrame(testCase, columns = ['horsepower', 'curb-weight', 'engine-size', 'highway-mpg', 'normalized-losses', 'symboling'])
SellingPrice = BestRR.predict(testDf)
# In Indian Currency (1 Dollar = 74.19 Rupees).
SellingPrice = 74.19
print("Selling Price of Your Friend's Car is {r,2f} Rupees.".format(SellingPrice[0]))
Selling Price of Your Friend's Car is 979,056.43 Rupees.
```