

## Linear Regression With Python.

### What is Regression:-

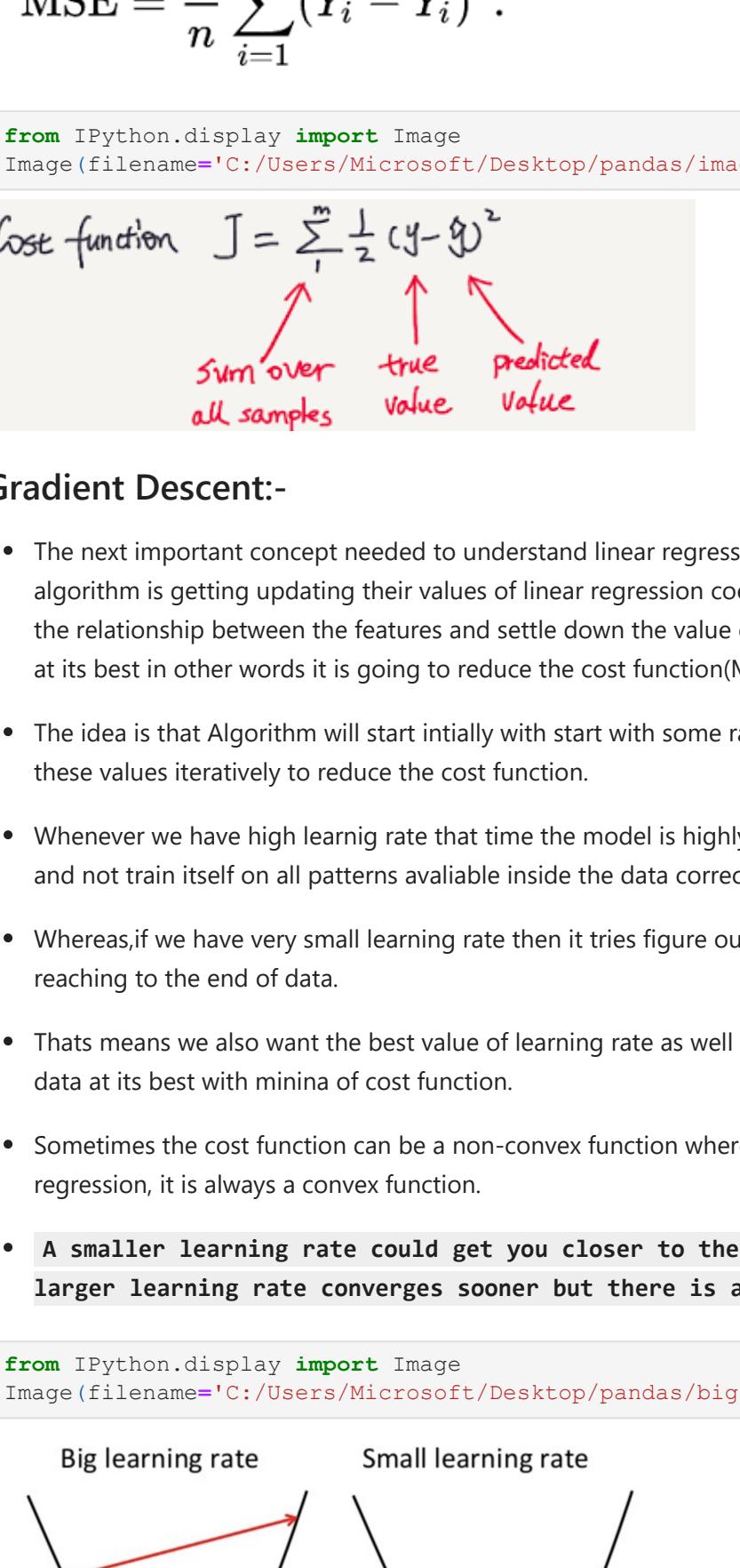
- Regression is a method of modelling a target value based on independent predictors. This method is mostly used for forecasting and finding out cause and effect relationship between variables.
- Regression techniques mostly differ based on the number of independent variables and the type of relationship between the independent and dependent variables.

### Linear Regression :-

- Simple linear regression is a type of regression analysis where the number of independent variables is one and there is a linear relationship between the independent(x) and dependent(y) variable. The red line in the below graph is referred to as the best fit straight line. Based on the given data points, we try to plot a line that gives us best fit line with least cost function or else we can say an error between the actual and predicted values. The line can be modelled based on the linear equation shown below.
- The motive of the linear regression algorithm is to find the optimal or best values for  $\beta_0$  and  $\beta_1$ .
- Before moving on to the algorithm, let's have a look at two important concepts you must know to better understand linear regression.

```
In [6]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/line.png',height=500,width=500)
```

Out[6]:



```
In [7]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/equation.png',height=500,width=500)
```

Out[7]:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Annotations pointing to parts of the equation:

- Dependent Variable:  $Y_i$
- Population Y intercept:  $\beta_0$
- Population Slope Coefficient:  $\beta_1$
- Independent Variable:  $X_i$
- Random Error term:  $\epsilon_i$
- Linear component:  $\beta_0 + \beta_1 X_i$
- Random Error component:  $\epsilon_i$

### Cost Function:-

- The cost function helps us to figure out the best possible values for  $\beta_0$  and  $\beta_1$  which would provide the best fit line for the data points. Since we want the best values for  $\beta_0$  and  $\beta_1$ .
- Algorithm initially needs to convert this randomly allocated values of  $\beta_0$  and  $\beta_1$  into their minima. It means we want the best values of these coefficient at optimal value where our model will give us at its best. That means we will reduce the error between the predicted and actual data points.
- The difference between the predicted values and ground truth measures the error difference. We square the error difference and sum over all data points and divide that value by the total number of data points.
- This provides the average squared error over all the data points. Therefore, this cost function is also known as the Mean Squared Error(MSE) function. Now, using this MSE function we are going to change the values of  $\beta_0$  and  $\beta_1$  such that the MSE value settles at its minima.

```
In [10]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/images.png',height=300,width=300)
```

Out[10]:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2.$$

```
In [14]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/images (1).png',height=400,width=400)
```

Out[14]:

$$\text{Cost function } J = \frac{1}{2} \sum_{i=1}^n (y - \hat{y})^2$$

Annotations:

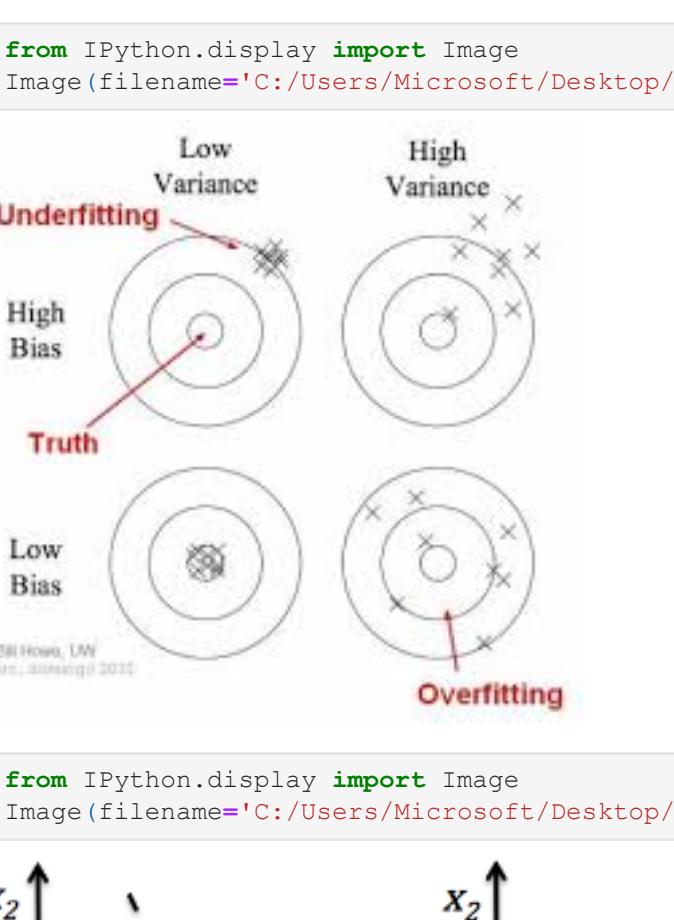
- Sum over all samples
- true value
- predicted value

### Gradient Descent:-

- The next important concept needed to understand linear regression is gradient descent. Gradient descent is a method where algorithm is getting updating their values of linear regression coefficient  $\beta_0$  and  $\beta_1$  gradually with certain learning rate in order to get the relationship between the features and settle down the value of regression coefficient at its best value where our model will perform at its best in other words it is going to reduce the cost function(MSE).
- The idea is that Algorithm will start initially with some random values of  $\beta_0, \beta_1, \dots, \beta_n$  and so on and then algorithm changes these values iteratively to reduce the cost function.
- Whenever we have high learning rate that time the model is highly unreliable in terms of prediction because it learned in zig-zag way and not train itself on all patterns available inside the data correctly.
- Whereas, if we have very small learning rate then it tries to figure out each pattern available data but it will stop learning soon before reaching to the end of data.
- That means we also want the best value of learning rate as well because it is going to help for understanding the pattern inside the data at its best with minima of cost function.
- Sometimes the cost function can be a non-convex function where you could settle at a local minima but in case of linear regression, it is always a convex function.
- A smaller learning rate could get you closer to the minima but takes more time to reach the minima, a larger learning rate converges sooner but there is a chance that you could overshoot the minima.

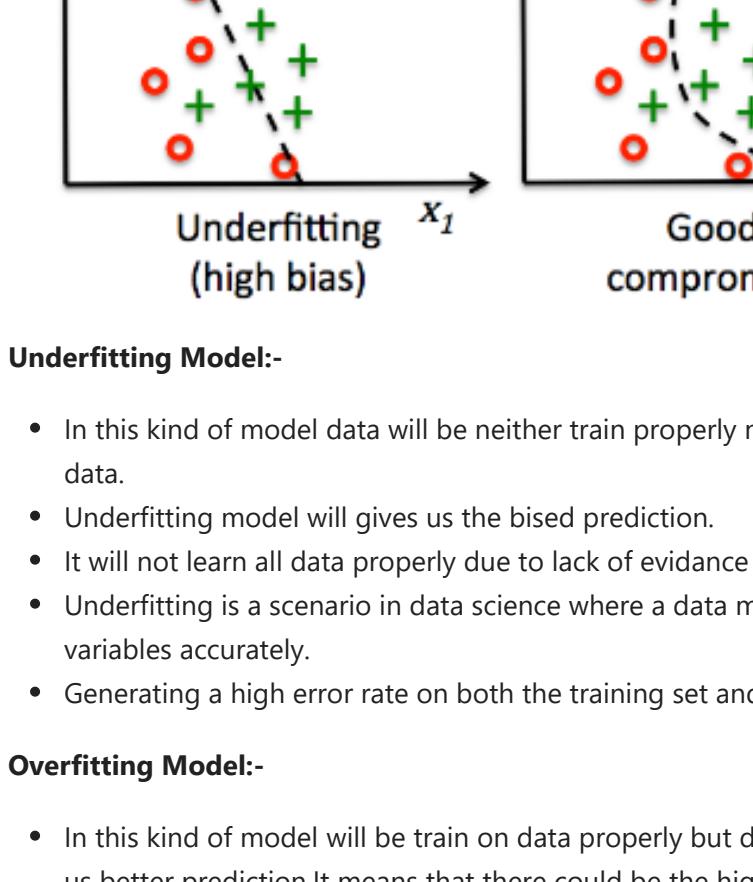
```
In [15]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/big.png',height=400,width=400)
```

Out[15]:



```
In [16]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/ball.png',height=400,width=400)
```

Out[16]:



### Code :-

- We use pandas library to read the train and test files. We retrieve the independent(x) and dependent(y) variables and since we have only one feature(x) we reshape them so that we could feed them into our linear regression model.
- We use scikit learn to import the LinearRegression model, we fit the model on the training data and predict the values for the testing data. We use R2 score to measure the accuracy of our model.

### Visualized our data:-

- By using the pairplot we can visualize our data to get the perfect behaviour.

### Example

#### 1. House Price Prediction on Boston Housing Datasets.

- You are the real estate agent and you want to predict the house price. It would be great if you can make some kind of automated system which can predict the price of the house based on various input which is known as features.
- Supervised Machine Learning Algorithm needs some data to train the model before making any prediction.
- We have a Boston Datasets.

### Where can The Linear Regression be used ?

- It is a very powerful technique and can be used to understand facts that can understand the probability. It can be used to predict forecast sales in the upcoming months by analyzing the sales data of the previous months. It can be used to gain various insights about the customer behaviour.

Here we are having the values of the  $x$  which are the independent to each other and  $y$  is the dependent variable.

$y$  used to calculate with the help of  $x$ .

Here  $x$  (Independent variable) and  $y$  (dependent variable) should always be in the form of continuous data in order to predict.

So, that we can say  $y$  is a function of  $x$ .

$y = f(x)$

For example:- If we have the housing price prediction, in that case the price of the house is on the  $y$ -axis, dependent variable, and rooms, blocks, halls and other things that would be counted in  $x$ -axis, that is independent variable.

### Types Of Linear Regression

- There are two kinds of regression
  - Simple linear regression.
  - Multiple regression.

### Simple linear regression

- Let's suppose if we go with the single input linear regression.
- In that we have the single input on the  $x$ -axis and according to this input we will get the straight line and it will give us the predicted value of house.

Once, the line will fit in the linear regression then we will find what is the goodness of the fit.

That means how well your model will behave on unseen data.

It behaves good or it produces more accuracy or good MSE reduction but it might not behave good on the unseen data.

So, we define goodness of fit on the unseen data and our goal is to minimize the total error and the total square error is a mean square error or mean absolute error.

That can be minimized by using the Gradient Descent.

### What gradient descent do in that ?

- We start fitting our regression line simultaneously with our data and we will calculate the mean squared error with this fitted line.

If error is in high amount then it will look forward again for getting the best position of itself inside the data where it can get the minimum error between predicted and actual values.

It is happening until and unless we will not get the best possible regression line or best fit line where cost function is at its minima.

### About the Line.

```
In [19]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/www.png',height=400,width=800)
```

Out[19]:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Annotations pointing to parts of the equation:

- Dependent Variable:  $Y_i$
- Population Y intercept:  $\beta_0$
- Population Slope Coefficient:  $\beta_1$
- Independent Variable:  $X_i$
- Random Error term:  $\epsilon_i$
- Linear component:  $\beta_0 + \beta_1 X_i$
- Random Error component:  $\epsilon_i$

### Multiple Linear Regression.

Multiple linear regression is the regression where we will use multiple independent inputs. Apart from that, here we will use the polynomial linear regression.

```
In [20]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/multiple_linear2.png',height=400,width=800)
```

Out[20]:

Simple Linear Regression:  $y = b_0 + b_1 * x_1$

Multiple Linear Regression:  $y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$

Annotations:

- Dependent variable (DV):  $y$
- Independent variables (IVs):  $x_1, x_2, \dots, x_n$
- Constant:  $b_0$
- Coefficients:  $b_1, b_2, \dots, b_n$

### Plane Fitting On the Linear Regression.

- If we use the multiple regression input, in that case it will use plane fitting instead of line fitting. Apart from that, here we will use the polynomial linear regression.

That will become multiple linear regression.

### Assessing the Performance of the model.

- Algorithm calculates the vertical distance from the regression line to the data points and determines which data points are close or far from the line that it fits.
- The difference between them is called residual.
- The line for which the error between the predicted values and observed values is minimum is called the best fit line of the regression line.
- This error is called residual error.

```
In [24]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/linearregression_re.jpg',height=400,width=500)
```

Out[24]:



### How do we determine Best Fit Line ?

```
In [37]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/ver.png',height=700,width=700)
```

Out[37]:



### Code :-

- We use pandas library to read the train and test files. We retrieve the independent(x) and dependent(y) variables and since we have only one feature(x) we reshape them so that we could feed them into our linear regression model.

- We use scikit learn to import the LinearRegression model, we fit the model on the training data and predict the values for the testing data. We use R2 score to measure the accuracy of our model.

### Visualized our data:-

- By using the pairplot we can visualize our data to get the perfect behavior.

### Example

#### 1. House Price Prediction on Boston Housing Datasets.

- You are the real estate agent and you want to predict the house price. It would be great if you can make some kind of automated system which can predict the price of the house based on various input which is known as features.

Supervised Machine Learning Algorithm needs some data to train the model before making any prediction.

We have a Boston Datasets.

### Where can The Linear Regression be used ?

- It is a very powerful technique and can be used to understand facts that can understand the probability. It can be used to predict forecast sales in the upcoming months by analyzing the sales data of the previous months. It can be used to gain various insights about the customer behaviour.

Here we are having the values of the  $x$  which are the independent to each other and  $y$  is the dependent variable.

$y$  used to calculate with the help of  $x$ .

Here  $x$  (Independent variable) and  $y$  (dependent variable) should always be in the form of continuous data in order to predict.

So, that we can say  $y$  is a function of  $x$ .

$y = f(x)$

For example:- If we have the housing price prediction, in that case the price of the house is on the  $y$ -axis, dependent variable, and rooms, blocks, halls and other things that would be counted in  $x$ -axis, that is independent variable.

### Types Of Linear Regression

- There are two kinds of regression
  - Simple linear regression.
  - Multiple regression.

### Simple linear regression

- Let's suppose if we go with the single input linear regression.

In that we have the single input on the  $x$ -axis and according to this input we will get the straight line and it will give us the predicted value of house.

Once, the line will fit in the linear regression then we will find what is the goodness of the fit.

That means how well your model will behave on unseen data.

It behaves good or it produces more accuracy or good MSE reduction but it might not behave good on the unseen data.

So, we define goodness of fit on the unseen data and our goal is to minimize the total error and the total square error is a mean square error or mean absolute error.

That can be minimized by using the Gradient Descent.

### What gradient descent do in that ?

- We start fitting our regression line simultaneously with our data and we will calculate the mean squared error with this fitted line.

If error is in high amount then it will look forward again for getting the best position of itself inside the data where it can get the minimum error between predicted and actual values.

It is happening until and unless we will not get the best possible regression line or best fit line where cost function is at its minima.

### About the Line.

```
In [19]: from IPython.display import Image  
Image(filename='C:/Users/Microsoft/Desktop/pandas/www.png',height=400,width=800)
```

Out[19]:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Annotations pointing to parts of the equation:

- Dependent Variable:  $Y_i$
- Population Y intercept:  $\beta_0$
- Population Slope Coefficient:  $\beta_1$
- Independent Variable:  $X_i$
- Random Error term:  $\epsilon_i$
- Linear component:  $\beta_0 + \beta_1 X_i$
- Random Error component:  $\epsilon_i$

### Multiple Linear Regression.

&lt;

## Introduction & Summary of Linear Regression.

- **Linear Regression**
- Regression means prediction dependent on the basis of independent variable.
- In the linear the value should be continuous format.
- In that will use the Euclidean distance to measure the error (error is the perpendicular distance from the base fit line or Regression Line)
- It will minimized by using the StandardScaler or gradient descent algorithm to find global minima there we will get the least error or root mean squared error.
- **Kinds Of Regression.**
- There are two types of LinearRegression

  - 1-Simple Regression
  - 2-Multiple regression i.e. Polynomial Regression.

- Both uses the linear relationship by using the straight line equation  $y = mx + c$
- **Bias-Variance Trade Off**
- Bias - It makes the prediction away from actual with low variance.
- Variance - It is dispersion of the actual datapoint away from the mean of that datapoint.
- Bias and Variance are mutually inversely proportional to each other.
- If Bias is high then Variance is low that time model complexity got decreases but due to this reason model undefined that it either poorly fit on our data.
- When the bias is high and variance is low that time model complexity got increases but due to this reason model undefined that model become Overfitted. It means that model train on each and every datapoint of the datasets but failed to generalize the model for new data.
- **Underfit and Overfit**
- Underfit - Whenever the variance is low and bias is high that means it is Underfitted.
- Overfitted - Whenever the variance is high but the bias is low that is Overfitted.
- **Identification Of Overfit and Underfit**
- Whenver the R2 score of training underfit model is greater than R2 score of testing based fit model that time we can say it is overfitted model or vice versa.

• **Errors**

- Mean Squared Error
- Mean Absolute error

### • Measure of Performance Count of the Linear Regression

• **R2score**

### • Assumptions Of linear Regression

• 1.Linearity - Graphical method regplot between the residue and prediction, 2nd Linear\_rainbow method 3rd is the mean the residual is closest to zero.

• 2.Normality - Jaque bera test (pvalue=5.99<5) and Graphical method by using the distribution plot of residual. Here null hypothesis is 5.99.

• 3.Homoscedasticity - Graphical method lmplot between the residue and prediction and Goldfread test or Bust wagon test. Here the value of the null hypothesis is 0.2.

• 4.No Autocorrelation- durbin-watson-test (0-2(positive Auto-correlation) and 2 to 4 (Negative Auto-correlation))

• 5.No Multicollinearity- Variance influence factor.

### • Feature Selection Done On The Basis Of High Correlation of the Independent with Dependedent Variable

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
warnings.filterwarnings('ignore')

In [2]: data=pd.read_csv('50_Startups.csv')

In [3]: data.head()

Out[3]:   R&D Spend    Adminstration Marketing Spend    State    Profit
0 165349.20 136897.80 471784.10 New York 192261.83
1 162597.70 151377.59 443986.53 California 191792.06
2 153441.51 101145.55 407934.54 Florida 191050.39
3 144372.41 118671.85 383199.62 New York 182901.99
4 142107.34 91391.77 366168.42 Florida 166187.94
```

```
In [4]: data.tail()
```

```
Out[4]:   R&D Spend    Adminstration Marketing Spend    State    Profit
45 1000.23 124153.04 1903.93 New York 64926.68
46 1315.46 115816.21 297114.46 Florida 49490.75
47 0.00 135426.92 0.00 California 42559.73
48 542.05 51743.15 0.00 New York 35673.41
49 0.00 116983.80 45173.06 California 14681.40
```

```
In [5]: data.rename(columns={'R&D Spend':'R&D','Marketing Spend':'Marketing'},inplace=True)
```

```
In [6]: data.head()
```

```
Out[6]:   R&D Administration Marketing State    Profit
0 165349.20 136897.80 471784.10 New York 192261.83
1 162597.70 151377.59 443986.53 California 191792.06
2 153441.51 101145.55 407934.54 Florida 191050.39
3 144372.41 118671.85 383199.62 New York 182901.99
4 142107.34 91391.77 366168.42 Florida 166187.94
```

```
In [7]: data.info()
RangeIndex: 50 entries, 0 to 49
Data columns (total 5 columns):
 #   Column  Non-Null Count  Dtype  
 0   R&D      50 non-null    float64 
 1   Administration 50 non-null    float64 
 2   Marketing  50 non-null    float64 
 3   State     50 non-null    object  
 4   Profit    50 non-null    float64 
dtypes: float64(4), object(1)
memory usage: 1.8+ KB
```

```
In [8]: data.describe()
```

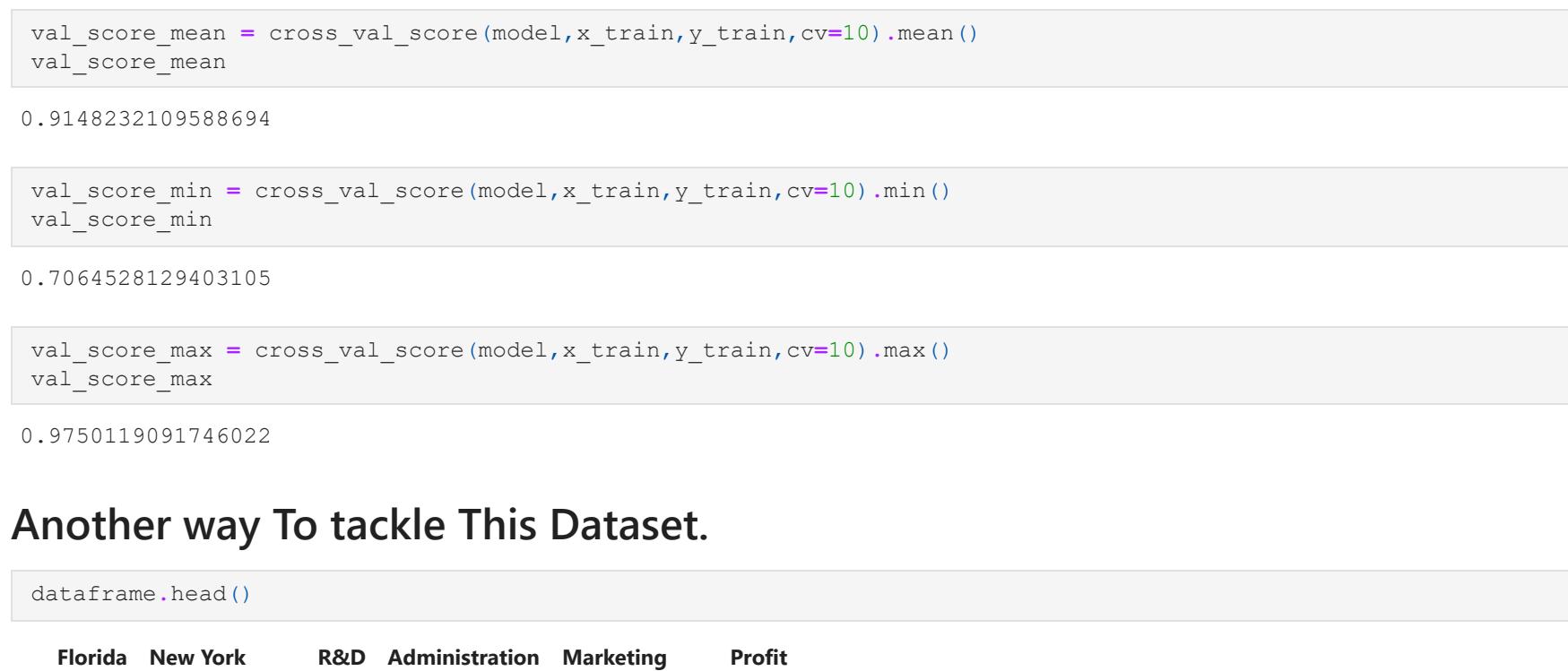
```
Out[8]:    R&D  Administration  Marketing  Profit
count  50.000000  50.000000  50.000000  50.000000
mean   737216.15602 121344.639600 2110259.7800 112012.639000
std    45902.54682 28017.802755 122290.310726 40306.18038
min    39963.70000 103730.875000 129300.132500 90138.902500
25%    73016.80000 122699.795000 212716.240000 10778.190000
50%    101660.80000 144842.180000 259469.085000 139765.977500
75%    165349.20000 182645.560000 471784.100000 192261.830000
max    165349.20000 182645.560000 471784.100000 192261.830000
```

```
In [9]: data.isnull().sum()
```

```
Out[9]: R&D    Administration    Marketing    State    Profit
Administration 0
Marketing 0
State 0
Profit 0
dtype: int64
```

### EDA

```
In [10]: data.plot()
plt.show()
```



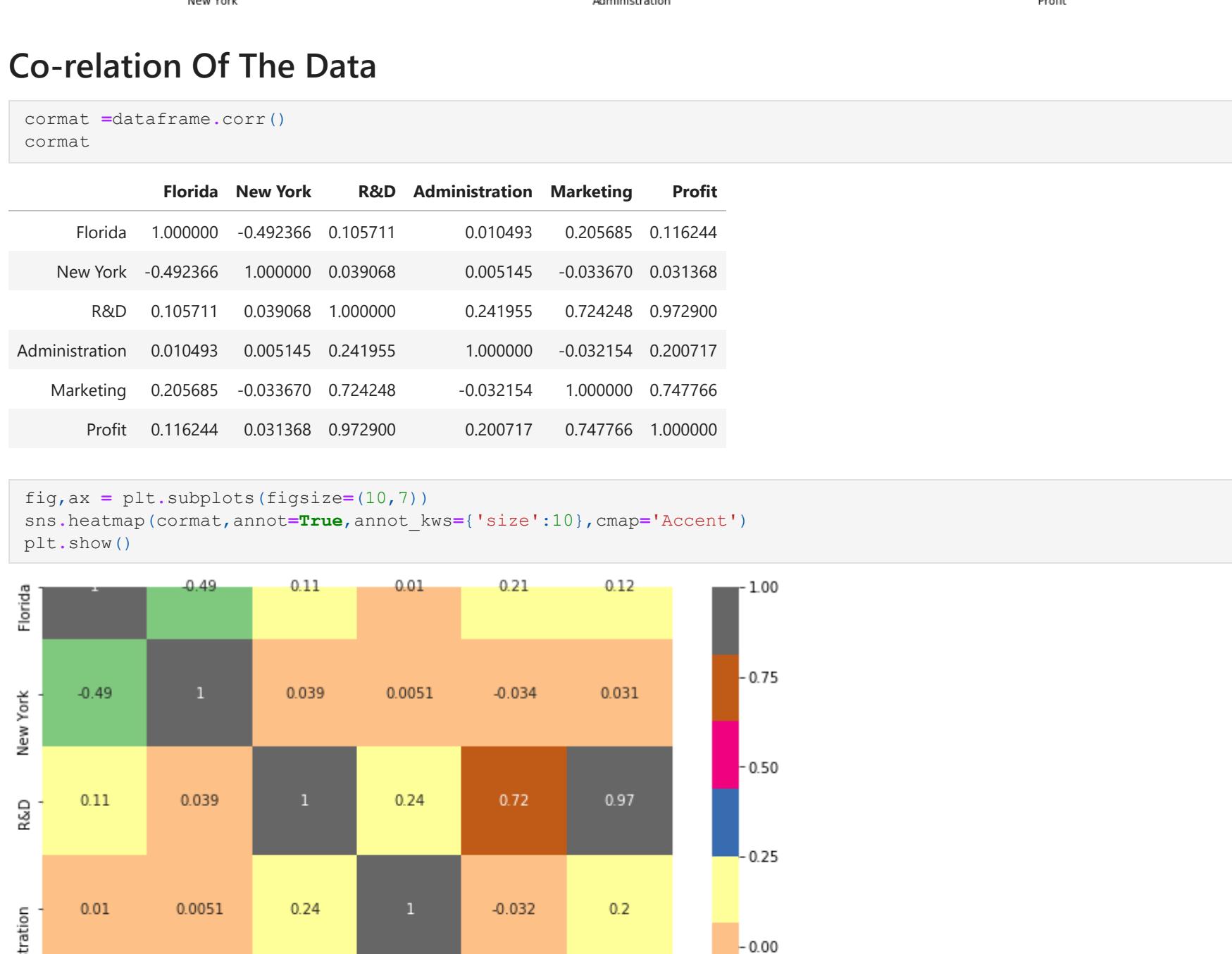
```
In [11]: fig,ax1,ax2,ax3,ax4=plt.subplots(nrows=4,ncols=1,figsize=(20,6))

ax1.plot(data['R&D'])
ax1.set_title('R&D')
ax2.plot(data['Profit'])
ax2.set_title('Profit')
ax3.plot(data['Administration'])
ax3.set_title('Administration')
ax4.plot(data['Marketing'])
ax4.set_title('Marketing')
plt.show()
```



```
In [12]: fig,(ax1,ax2,ax3)=plt.subplots(ncols=3,nrows=1,figsize=(15,5))

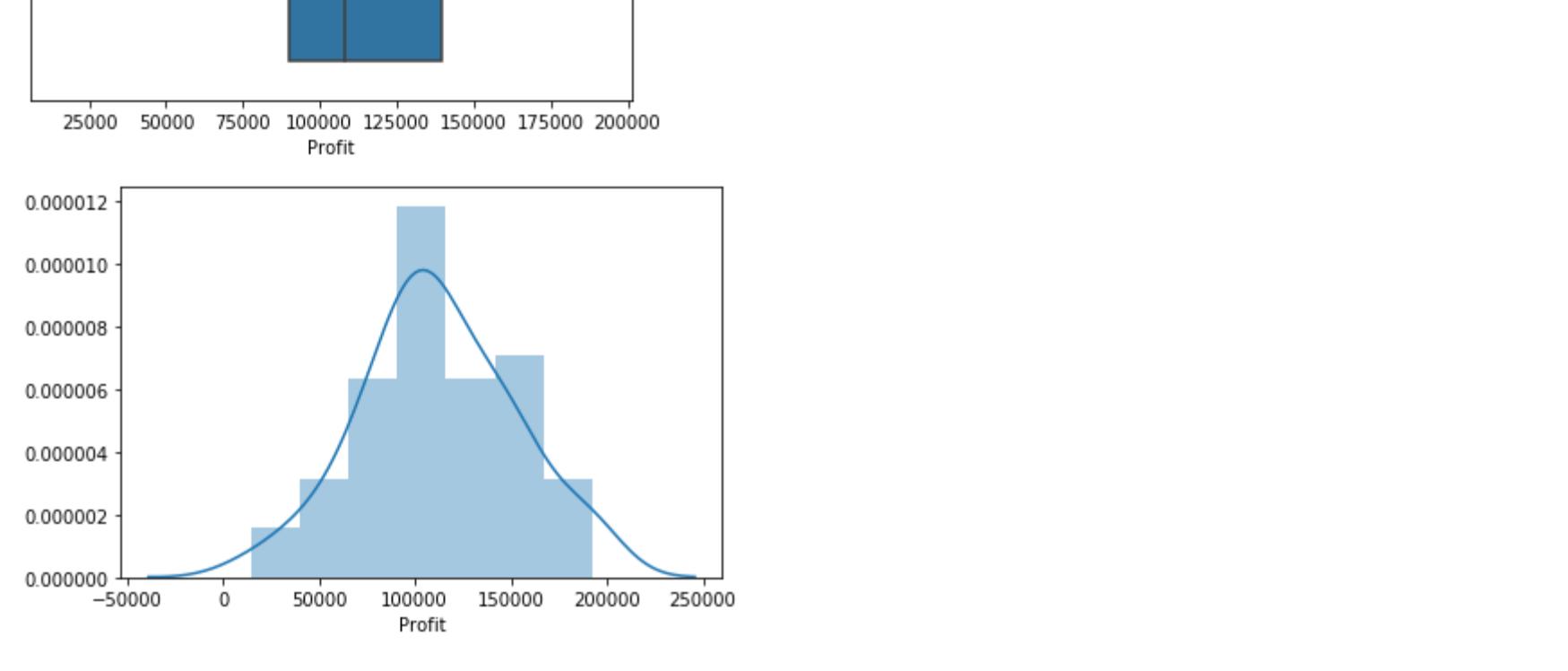
ax1.plot(data['R&D'].sort_values(ascending=False),data['Profit'].sort_values(ascending=False))
ax1.set_title('R&D Vs Profit')
plt.xticks(rotation=90)
ax2.plot(data['Marketing'].sort_values(ascending=False),data['Profit'].sort_values(ascending=False))
ax2.set_title('Marketing Vs Profit')
plt.xticks(rotation=90)
ax3.plot(data['Administration'].sort_values(ascending=False),data['Profit'].sort_values(ascending=False))
ax3.set_title('Administration vs Profit')
plt.xticks(rotation=90)
plt.show()
```



```
In [13]: data['State'].value_counts()
```

```
Out[13]: California    17
New York     17
Florida      16
Name: State, dtype: int64
```

```
In [14]: sns.countplot(data['State'])
plt.show()
```



```
In [15]: data.groupby(data['State']).count()
```

```
Out[15]:   R&D  Administration  Marketing  Profit
State
California 17
New York 17
Florida 16
```

```
In [16]: dig = pd.concat([data['State'],pd.get_dummies(data['State'])],axis=1)
dig.drop(['State'],axis=1,inplace=True)
dig.head()
```

```
In [17]: df = dig[['R&D','Administration','Marketing','State','Profit']]
df['New York']=0
df['California']=1
df['Florida']=2
```

```
In [18]: sns.catplot(y='Administration',x='Profit',data=df,hue='State',height=8)
plt.show()
```



```
In [19]: sns.jointplot(y='Administration',x='Profit',data=df,color='g',kind='kde')
plt.show()
```



```
In [20]: plt.figure(figsize=(20,8),dpi=100)
sns.barplot('R&D','Administration',hue='State')
plt.title('Administration vs Profit')
plt.show()
plt.tight_layout()
```

```
In [21]: plt.figure(figsize=(20,8),dpi=100)
sns.barplot('R&D','Marketing',hue='State')
plt.title('R&D Vs Profit')
plt.show()
plt.tight_layout()
```

```
In [22]: plt.figure(figsize=(20,8),dpi=100)
sns.barplot('R&D','Profit',hue='State')
plt.title('R&D Vs Profit')
plt.show()
plt.tight_layout()
```

```
In [23]: Florida New York R&D Administration Marketing Profit
0 0 1 165349.20 136897.80 471784.10 192261.83
1 0 0 162597.70 151377.59 443986.53 191792.06
2 1 0 153441.51 101145.55 407934.54 191050.39
3 0 1 144372.41 118671.85 383199.62 182901.99
4 1 0 142107.34 91391.77 366168.42 166187.94
```

### Start To Prepare Model

```
In [24]: from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import StandardScaler
```

```
In [25]: x=df.drop(['Profit'],axis=1)
y=df['Profit']
```

```
In [26]: x=np.array(x)
y=np.array(y)
```

```
In [27]: Train_test_split
```

```
In [28]: x_train,y_train,x_test,y_test = train_test_split(x,y,train_size=0.8,random_state=42)
```

```
In [29]: x=x_train
y=y_train
```

```
In [30]: x=x.reshape(-1,1)
y=y.reshape(-1,1)
```

```
In [31]: x=x.reshape(-1,1)
y=y.reshape(-1,1)
```

```
In [32]: x=x.reshape(-1,1)
y=y.reshape(-1,1)
```

```
In [33]: model=LinearRegression()
model.fit(x,y)
```

```
In [34]: print("Slope Of The Line :-",model.coef_)
print("Intercept of the line :-",model.intercept_)
```

```
In [35]: Slope Of The Line :- 4.4777572e+02
Intercept of the line :- 2.0178910e+04
3.3861758e+03
```

```
In [36]: print("Accuracy over the training data :-",model.score(x_train,y_train))
Accuracy over the training data :- 0.953701993240556
```

```
In [37]: print("Accuracy over the training data :-",model.score(x_test,y_test))
Accuracy over the training data :- 0.998726614319837
```

```
In [38]: from sklearn.model_selection import cross_val_score
val_score_mean = cross_val_score(model,x_train,y_train,cv=10).mean()
```

```
In [39]: val_score_min = cross_val_score(model,x_train,y_train,cv=10).min()
```

```
In [40]: val_score_max = cross_val_score(model,x_train,y_train,cv=10).max()
```

```
In [41]: 0.7064528129403105
```

```
In [42]: 0.975109691746022
```

### Another Way To Tackle This Dataset.

```
In [43]: dataframe.head()
```

```
Out[43]:   Florida New York R&D Administration Marketing Profit
0 0 1 165349.20 136897.80 471784.10 192261.83
1 0 0 162597.70 151377.59 443986.53 191792.06
2 1 0 153441.51 101145.55 407934.54 191050.39
3 0 1 144372.41 118671.85 383199.62 182901.99
4 1 0 142107.34 91391.77 366168.42 166187.94
```

```
In [44]: corr=df.corr()
```

```
Out[44]:   Florida New York R&D Administration Marketing Profit
Florida 1.000000 0.492366 0.105711 0.205685 0.116244
New York 0.492366 1.000000 0.039068 0.005145 -0.033670 0.031368
R&D 0.104911 0.039068 1.000000 0.241955 0.127424 0.972900
Administration 0.005145 0.241955 1.000000 -0.032154 0.190084 0.200717
Marketing 0.205685 0.116244 0.241955 1.000000 0.747766 1.000000
Profit 0.116244 0.031368 0.127424 0.747766 1.000000 0.000000
```

```
In [45]: dig=plt.subplots(figsize=(10,7))
dig.heatmap(df.corr(),cmap='Accent')
plt.show()
```



```
In [46]: array(['Florida', 'New York', 'R&D', 'Administration', 'Marketing', 'Profit'],
      dtype='object')
```

```
In [47]: sns.boxplot(dataframe['Profit'])
sns.distplot(dataframe['Profit'])
plt.show()
```



```
In [48]: dataframe['Administration'].hist()
plt.show()
```

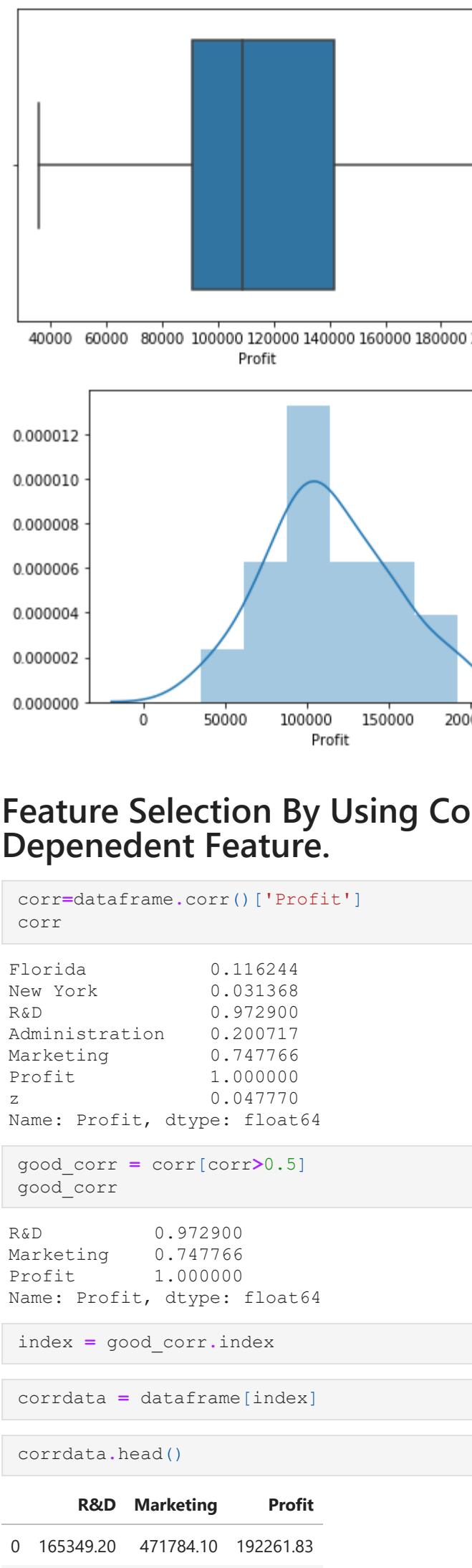


```
In [49]: dataframe['Marketing'].hist()
plt.show()
```



```
In [50]: dataframe['R&D'].hist()
plt.show()
```





## Feature Selection By Using Correlation Between the Independent and Dependent Feature.

```
In [61]: corr=corrdato.corr()['Profit']
```

```
Out[61]:
```

Florida	0.116244
New York	0.151938
R&D	0.372300
Administration	0.200717
Marketing	0.147766
Profit	1.000000
z	0.047770

```
Name: Profit, dtype: float64
```

```
In [62]: good_corr = corr[corr>0.5]
```

```
Out[62]:
```

R&D	0.97900
Marketing	0.747766
Profit	1.000000
z	0.047770

```
Name: Profit, dtype: float64
```

```
In [63]: index = good_corr.index
```

```
In [64]: corrdato = corrdato[index]
```

```
In [65]: corrdato.head()
```

```
Out[65]:
```

	R&D	Marketing	Profit
0	163349.20	471784.10	192261.83
1	162597.70	443988.53	191792.06
2	153441.51	407934.54	191050.39
3	144372.41	383199.62	182901.99
4	142107.34	366168.42	166187.94

```
In [66]: plt.figure(figsize=(14,8))
sns.pairplot(corrdato)
plt.show()
```

<Figure size 1008x576 with 0 Axes>

```
In [67]: plt.figure(figsize=(10,6))
sns.heatmap(corrdato.corr(), annot=True)
plt.show()
```

```
Out[67]:
```

R&D	1	-0.97	0.047770
Marketing	-0.97	1	0.75
Profit	0.047770	1	1.00

**Create model**

```
In [68]: X = corrdato.drop(columns=['Profit'])
Y = corrdato['Profit']
```

```
In [69]: X_train,X_test,Y_train,Y_test = train_test_split(X,Y,train_size=0.80,random_state=42)
```

```
In [70]: (X_train.shape,Y_train.shape)
```

```
Out[70]: ((40, 3), (40,))
```

```
In [71]: scalar = StandardScaler()
X_train=scalar.fit_transform(X_train)
X_test = scalar.transform(X_test)
```

**R2 score**

R2 score is nothing but the **1-(mean Absolute Square)**.

- If the r2 score is nearer to one then we can say the model is good
- If the r2 score is nearer to zero then we can say the model is not so much good.

**What is R-Squared?**

R-squared (R<sup>2</sup>) is a statistical measure that represents the proportion of the variance for a dependent variable that's explained by an independent variable or variables in a regression model. Whereas correlation explains the strength of the relationship between an independent and dependent variable, R-squared explains to what extent the variance of one variable explains the variance of the second variable. So, if the R<sup>2</sup> of a model is 0.50, then approximately half of the observed variation can be explained by the model's inputs.

```
In [72]: from IPython.display import Image
image=filename='C:/Users/Microsoft/Desktop/pandas/download (2).png'
```

Out[72]:

**Sum Squared Regression Error**

$$R^2 = 1 - \frac{SS_{\text{Regression}}}{SS_{\text{Total}}}$$

Sum Squared Total Error

**Below we can see the**

- y= Actual Point in data
- y<sub>pred</sub> = Predicted Y in the data
- y<sub>mean</sub> = Mean of the actual data.

```
In [73]: from IPython.display import Image
image=filename='C:/Users/Microsoft/Desktop/pandas/Data+Analysis.jpg'
```

Out[73]:

## Data Analysis

- R<sup>2</sup> or Coefficient of Determination. Equals the proportion of the variance in the dependent variable Y that is explained through the relationship with the Independent variable X.

- Explained Variance = Total Variance - Unexplained Variance. We state this as a proportion:

$$\text{Unadjusted } R^2 = 1 - \frac{\sum(Y_i - \hat{Y}_i)^2}{\sum(Y_i - \bar{Y})^2} \quad \text{Adjusted: } R^2 = 1 - \frac{S_{\text{res}}^2}{S_y^2}$$

Adjusted R<sup>2</sup> - adjusted for complexity by the degrees of freedom. Unadjusted R<sup>2</sup> becomes larger as more variables are added to the equation (decreases the sum of errors in the denominator). The use of an unadjusted R<sup>2</sup> may result in believing that additional variables are useful when they are not.

```
In [74]: from IPython.display import Image
image=filename='C:/Users/Microsoft/Desktop/pandas/regression-analysis-using-python-47-638.jpg'
```

Out[74]:

### Assessing goodness of fit: R<sup>2</sup>

For multiple linear regression, the coefficient of determination R<sup>2</sup> is calculated as

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

where

- $SS_{\text{res}} = \sum(Y_i - \hat{Y}_i)^2$  is the sum of the square of the residuals
- $SS_{\text{tot}} = \sum(Y_i - \bar{Y})^2$  is the total sum of squares

$y_i$  are the observations, for  $i = 1 \dots n$

$\hat{Y}_i$  are the predictions, for  $i = 1 \dots n$

$\bar{Y}$  is the mean of the observations

The better the fit, the closer R<sup>2</sup> is to 1

R<sup>2</sup> measures the proportion of variance in the observed data that is explained by the model

Area of red squares: squared residuals with respect to the average value

Area of blue squares: squared residuals with respect to the linear regression

RISK  
ENGINEERING

```
In [75]: model = LinearRegression()
model.fit(X_train,Y_train)
y_pred = model.predict(X_test)
print('Accuracy :',r2_score(Y_test,y_pred))
```

```
Out[75]:
```

Accuracy :	0.9168381183550247
------------	--------------------

```
In [76]: df = pd.DataFrame({'Actual':Y_test,'Predicted':y_pred})
```

Out[76]:

### Matrix Measures Of The Linear Regression.

- R2 score.
- Mean Absolute Error(MAE)
- Mean Squared Error(MSE)

```
In [77]: from sklearn.metrics import r2_score,mean_squared_error,mean_absolute_error
```

```
In [78]: score = r2_score(Y_test,y_pred)
mae = mean_absolute_error(Y_test,y_pred)
mse = mean_squared_error(Y_test,y_pred)
```

```
Out[78]:
```

r2_score :	0.9168381183550247
Mean absolute error :	6459.181714685925
Mean squared Error :	67343832.5897961

```
In [79]: MSE = np.sqrt(mse)
MSE
```

```
Out[79]:
```

8206.328813165854
-------------------

```
In [80]: col = 2
row = 1
fig,ax = plt.subplots(row,col,figsize=(15,7))

cols = corrdato.columns
index = range(len(cols)):
    sns.regplot(x=corrdato[cols[index]],y=corrdato['Profit'],x_order=index)
    index = index+1
plt.tight_layout()
```

Out[80]:

To identify the linearity

To identify the Normality

To identify the linearity

```
In [27]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
import warnings
warnings.filterwarnings('ignore')

In [28]: data = pd.read_csv('kc_house_train_data.csv')

In [29]: data.head()

Out[29]:   Id MSSubClass MSZoning LotFrontage LotArea Street Alley LotShape LandContour Utilities ... PoolArea PoolQC Fence MiscFeature
0  1     60       RL      65.0    8450  Pave  Reg    Reg  Lvl  AllPub   ...        0  NaN  NaN  NaN
1  2     60       RL      80.0    9600  Pave  NaN  Reg  Lvl  AllPub   ...        0  NaN  NaN  NaN
2  3     60       RL      68.0    9550  Pave  NaN  IR1  Lvl  AllPub   ...        0  NaN  NaN  NaN
3  4     70       RL      60.0    9550  Pave  NaN  IR1  Lvl  AllPub   ...        0  NaN  NaN  NaN
4  5     60       RL      84.0   14260  Pave  NaN  IR1  Lvl  AllPub   ...        0  NaN  NaN  NaN
5 rows × 81 columns

In [30]: data.tail()

Out[30]:   Id MSSubClass MSZoning LotFrontage LotArea Street Alley LotShape LandContour Utilities ... PoolArea PoolQC Fence MiscFeature
1455 1456     60       RL      62.0    7917  Pave  NaN  Reg  Lvl  AllPub   ...        0  NaN  NaN  NaN
1456 1457     20       RL      85.0   13175  Pave  NaN  Reg  Lvl  AllPub   ...        0  NaN  Mdpv
1457 1458     70       RL      60.0    9042  Pave  NaN  Reg  Lvl  AllPub   ...        0  NaN  Gdpv
1458 1459     20       RL      68.0    9717  Pave  NaN  IR1  Lvl  AllPub   ...        0  NaN  NaN  NaN
1459 1460     20       RL      75.0    9937  Pave  NaN  Reg  Lvl  AllPub   ...        0  NaN  NaN  NaN
5 rows × 81 columns

In [31]: data.columns

Out[31]: Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street', 'Alley', 'LotShape', 'LandContour', 'Utilities', '... ', 'PoolArea', 'PoolQC', 'Fence', 'MiscFeature'],
       dtype='object')

In [32]: len(data.columns)

Out[32]: 81

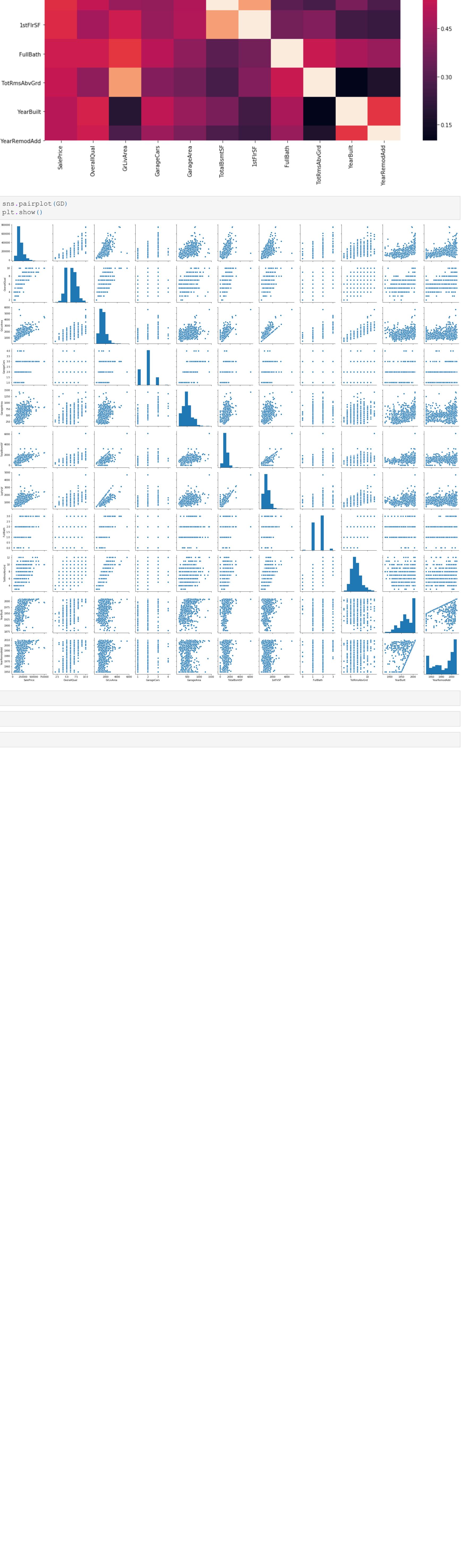
In [33]: data.shape

Out[33]: (1460, 81)

In [34]: print(data.info())

Out[34]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Id              1460 non-null   int64  
 1   MSSubClass      1460 non-null   object 
 2   MSZoning        1460 non-null   object 
 3   LotFrontage     1460 non-null   float64
 4   LotArea         1460 non-null   float64
 5   Street          1460 non-null   object 
 6   Alley           91 non-null    object 
 7   LotShape        1460 non-null   object 
 8   LandContour    1460 non-null   object 
 9   Utilities       1460 non-null   object 
 10  ...             ...             ...
 11  PoolArea        1460 non-null   float64
 12  PoolQC          1460 non-null   object 
 13  Fence           1460 non-null   object 
 14  MiscFeature    1460 non-null   object 
 15  SalePrice       1460 non-null   float64
 16  SaleType        1460 non-null   object 
 17  SaleCondition   1460 non-null   object 
 18  SaleDate        1460 non-null   datetime64[ns]
 19  OverallQual    1460 non-null   float64
 20  OverallCond    1460 non-null   float64
 21  YearBuilt       1460 non-null   int64  
 22  YearRemodAdd   1460 non-null   int64  
 23  MSVnArea        1460 non-null   float64
 24  BsmtFinSF1     1460 non-null   float64
 25  BsmtFinSF2     1460 non-null   float64
 26  BsmtUnfSF      1460 non-null   float64
 27  TotalBsmtSF    1460 non-null   float64
 28  BsmtFullBath   1460 non-null   float64
 29  BsmtHalfBath   1460 non-null   float64
 30  FullBath        1460 non-null   float64
 31  HalfBath        1460 non-null   float64
 32  Bedrms          1460 non-null   float64
 33  BsmtBedrms     1460 non-null   float64
 34  BsmtFinType1   1460 non-null   object 
 35  BsmtFinType2   1460 non-null   object 
 36  BsmtCond       1460 non-null   object 
 37  BsmtExposure   1460 non-null   object 
 38  BsmtExposure2  1460 non-null   object 
 39  BsmtFinPct    1460 non-null   float64
 40  BsmtUnfPct    1460 non-null   float64
 41  BsmtFinType1_2 1460 non-null   object 
 42  BsmtFinType2_2 1460 non-null   object 
 43  BsmtCond_2     1460 non-null   object 
 44  BsmtExposure_2 1460 non-null   object 
 45  BsmtExposure2_2 1460 non-null   object 
 46  BsmtFinPct_2  1460 non-null   float64
 47  BsmtUnfPct_2  1460 non-null   float64
 48  BsmtFinType1_3 1460 non-null   object 
 49  BsmtFinType2_3 1460 non-null   object 
 50  BsmtCond_3     1460 non-null   object 
 51  BsmtExposure_3 1460 non-null   object 
 52  BsmtExposure2_3 1460 non-null   object 
 53  BsmtFinPct_3  1460 non-null   float64
 54  BsmtUnfPct_3  1460 non-null   float64
 55  BsmtFinType1_4 1460 non-null   object 
 56  BsmtFinType2_4 1460 non-null   object 
 57  BsmtCond_4     1460 non-null   object 
 58  BsmtExposure_4 1460 non-null   object 
 59  BsmtExposure2_4 1460 non-null   object 
 60  BsmtFinPct_4  1460 non-null   float64
 61  BsmtUnfPct_4  1460 non-null   float64
 62  BsmtFinType1_5 1460 non-null   object 
 63  BsmtFinType2_5 1460 non-null   object 
 64  BsmtCond_5     1460 non-null   object 
 65  BsmtExposure_5 1460 non-null   object 
 66  BsmtExposure2_5 1460 non-null   object 
 67  BsmtFinPct_5  1460 non-null   float64
 68  BsmtUnfPct_5  1460 non-null   float64
 69  BsmtFinType1_6 1460 non-null   object 
 70  BsmtFinType2_6 1460 non-null   object 
 71  BsmtCond_6     1460 non-null   object 
 72  BsmtExposure_6 1460 non-null   object 
 73  BsmtExposure2_6 1460 non-null   object 
 74  BsmtFinPct_6  1460 non-null   float64
 75  BsmtUnfPct_6  1460 non-null   float64
 76  BsmtFinType1_7 1460 non-null   object 
 77  BsmtFinType2_7 1460 non-null   object 
 78  BsmtCond_7     1460 non-null   object 
 79  BsmtExposure_7 1460 non-null   object 
 80  BsmtExposure2_7 1460 non-null   object 
 81  BsmtFinPct_7  1460 non-null   float64
 82  BsmtUnfPct_7  1460 non-null   float64
 83  BsmtFinType1_8 1460 non-null   object 
 84  BsmtFinType2_8 1460 non-null   object 
 85  BsmtCond_8     1460 non-null   object 
 86  BsmtExposure_8 1460 non-null   object 
 87  BsmtExposure2_8 1460 non-null   object 
 88  BsmtFinPct_8  1460 non-null   float64
 89  BsmtUnfPct_8  1460 non-null   float64
 90  BsmtFinType1_9 1460 non-null   object 
 91  BsmtFinType2_9 1460 non-null   object 
 92  BsmtCond_9     1460 non-null   object 
 93  BsmtExposure_9 1460 non-null   object 
 94  BsmtExposure2_9 1460 non-null   object 
 95  BsmtFinPct_9  1460 non-null   float64
 96  BsmtUnfPct_9  1460 non-null   float64
 97  BsmtFinType1_10 1460 non-null   object 
 98  BsmtFinType2_10 1460 non-null   object 
 99  BsmtCond_10    1460 non-null   object 
 100  BsmtExposure_10 1460 non-null   object 
 101  BsmtExposure2_10 1460 non-null   object 
 102  BsmtFinPct_10 1460 non-null   float64
 103  BsmtUnfPct_10 1460 non-null   float64
 104  BsmtFinType1_11 1460 non-null   object 
 105  BsmtFinType2_11 1460 non-null   object 
 106  BsmtCond_11    1460 non-null   object 
 107  BsmtExposure_11 1460 non-null   object 
 108  BsmtExposure2_11 1460 non-null   object 
 109  BsmtFinPct_11 1460 non-null   float64
 110  BsmtUnfPct_11 1460 non-null   float64
 111  BsmtFinType1_12 1460 non-null   object 
 112  BsmtFinType2_12 1460 non-null   object 
 113  BsmtCond_12    1460 non-null   object 
 114  BsmtExposure_12 1460 non-null   object 
 115  BsmtExposure2_12 1460 non-null   object 
 116  BsmtFinPct_12 1460 non-null   float64
 117  BsmtUnfPct_12 1460 non-null   float64
 118  BsmtFinType1_13 1460 non-null   object 
 119  BsmtFinType2_13 1460 non-null   object 
 120  BsmtCond_13    1460 non-null   object 
 121  BsmtExposure_13 1460 non-null   object 
 122  BsmtExposure2_13 1460 non-null   object 
 123  BsmtFinPct_13 1460 non-null   float64
 124  BsmtUnfPct_13 1460 non-null   float64
 125  BsmtFinType1_14 1460 non-null   object 
 126  BsmtFinType2_14 1460 non-null   object 
 127  BsmtCond_14    1460 non-null   object 
 128  BsmtExposure_14 1460 non-null   object 
 129  BsmtExposure2_14 1460 non-null   object 
 130  BsmtFinPct_14 1460 non-null   float64
 131  BsmtUnfPct_14 1460 non-null   float64
 132  BsmtFinType1_15 1460 non-null   object 
 133  BsmtFinType2_15 1460 non-null   object 
 134  BsmtCond_15    1460 non-null   object 
 135  BsmtExposure_15 1460 non-null   object 
 136  BsmtExposure2_15 1460 non-null   object 
 137  BsmtFinPct_15 1460 non-null   float64
 138  BsmtUnfPct_15 1460 non-null   float64
 139  BsmtFinType1_16 1460 non-null   object 
 140  BsmtFinType2_16 1460 non-null   object 
 141  BsmtCond_16    1460 non-null   object 
 142  BsmtExposure_16 1460 non-null   object 
 143  BsmtExposure2_16 1460 non-null   object 
 144  BsmtFinPct_16 1460 non-null   float64
 145  BsmtUnfPct_16 1460 non-null   float64
 146  BsmtFinType1_17 1460 non-null   object 
 147  BsmtFinType2_17 1460 non-null   object 
 148  BsmtCond_17    1460 non-null   object 
 149  BsmtExposure_17 1460 non-null   object 
 150  BsmtExposure2_17 1460 non-null   object 
 151  BsmtFinPct_17 1460 non-null   float64
 152  BsmtUnfPct_17 1460 non-null   float64
 153  BsmtFinType1_18 1460 non-null   object 
 154  BsmtFinType2_18 1460 non-null   object 
 155  BsmtCond_18    1460 non-null   object 
 156  BsmtExposure_18 1460 non-null   object 
 157  BsmtExposure2_18 1460 non-null   object 
 158  BsmtFinPct_18 1460 non-null   float64
 159  BsmtUnfPct_18 1460 non-null   float64
 160  BsmtFinType1_19 1460 non-null   object 
 161  BsmtFinType2_19 1460 non-null   object 
 162  BsmtCond_19    1460 non-null   object 
 163  BsmtExposure_19 1460 non-null   object 
 164  BsmtExposure2_19 1460 non-null   object 
 165  BsmtFinPct_19 1460 non-null   float64
 166  BsmtUnfPct_19 1460 non-null   float64
 167  BsmtFinType1_20 1460 non-null   object 
 168  BsmtFinType2_20 1460 non-null   object 
 169  BsmtCond_20    1460 non-null   object 
 170  BsmtExposure_20 1460 non-null   object 
 171  BsmtExposure2_20 1460 non-null   object 
 172  BsmtFinPct_20 1460 non-null   float64
 173  BsmtUnfPct_20 1460 non-null   float64
 174  BsmtFinType1_21 1460 non-null   object 
 175  BsmtFinType2_21 1460 non-null   object 
 176  BsmtCond_21    1460 non-null   object 
 177  BsmtExposure_21 1460 non-null   object 
 178  BsmtExposure2_21 1460 non-null   object 
 179  BsmtFinPct_21 1460 non-null   float64
 180  BsmtUnfPct_21 1460 non-null   float64
 181  BsmtFinType1_22 1460 non-null   object 
 182  BsmtFinType2_22 1460 non-null   object 
 183  BsmtCond_22    1460 non-null   object 
 184  BsmtExposure_22 1460 non-null   object 
 185  BsmtExposure2_22 1460 non-null   object 
 186  BsmtFinPct_22 1460 non-null   float64
 187  BsmtUnfPct_22 1460 non-null   float64
 188  BsmtFinType1_23 1460 non-null   object 
 189  BsmtFinType2_23 1460 non-null   object 
 190  BsmtCond_23    1460 non-null   object 
 191  BsmtExposure_23 1460 non-null   object 
 192  BsmtExposure2_23 1460 non-null   object 
 193  BsmtFinPct_23 1460 non-null   float64
 194  BsmtUnfPct_23 1460 non-null   float64
 195  BsmtFinType1_24 1460 non-null   object 
 196  BsmtFinType2_24 1460 non-null   object 
 197  BsmtCond_24    1460 non-null   object 
 198  BsmtExposure_24 1460 non-null   object 
 199  BsmtExposure2_24 1460 non-null   object 
 200  BsmtFinPct_24 1460 non-null   float64
 201  BsmtUnfPct_24 1460 non-null   float64
 202  BsmtFinType1_25 1460 non-null   object 
 203  BsmtFinType2_25 1460 non-null   object 
 204  BsmtCond_25    1460 non-null   object 
 205  BsmtExposure_25 1460 non-null   object 
 206  BsmtExposure2_25 1460 non-null   object 
 207  BsmtFinPct_25 1460 non-null   float64
 208  BsmtUnfPct_25 1460 non-null   float64
 209  BsmtFinType1_26 1460 non-null   object 
 210  BsmtFinType2_26 1460 non-null   object 
 211  BsmtCond_26    1460 non-null   object 
 212  BsmtExposure_26 1460 non-null   object 
 213  BsmtExposure2_26 1460 non-null   object 
 214  BsmtFinPct_26 1460 non-null   float64
 215  BsmtUnfPct_26 1460 non-null   float64
 216  BsmtFinType1_27 1460 non-null   object 
 217  BsmtFinType2_27 1460 non-null   object 
 218  BsmtCond_27    1460 non-null   object 
 219  BsmtExposure_27 1460 non-null   object 
 220  BsmtExposure2_27 1460 non-null   object 
 221  BsmtFinPct_27 1460 non-null   float64
 222  BsmtUnfPct_27 1460 non-null   float64
 223  BsmtFinType1_28 1460 non-null   object 
 224  BsmtFinType2_28 1460 non-null   object 
 225  BsmtCond_28    1460 non-null   object 
 226  BsmtExposure_28 1460 non-null   object 
 227  BsmtExposure2_28 1460 non-null   object 
 228  BsmtFinPct_28 1460 non-null   float64
 229  BsmtUnfPct_28 1460 non-null   float64
 230  BsmtFinType1_29 1460 non-null   object 
 231  BsmtFinType2_29 1460 non-null   object 
 232  BsmtCond_29    1460 non-null   object 
 233  BsmtExposure_29 1460 non-null   object 
 234  BsmtExposure2_29 1460 non-null   object 
 235  BsmtFinPct_29 1460 non-null   float64
 236  BsmtUnfPct_29 1460 non-null   float64
 237  BsmtFinType1_30 1460 non-null   object 
 238  BsmtFinType2_30 1460 non-null   object 
 239  BsmtCond_30    1460 non-null   object 
 240  BsmtExposure_30 1460 non-null   object 
 241  BsmtExposure2_30 1460 non-null   object 
 242  BsmtFinPct_30 1460 non-null   float64
 243  BsmtUnfPct_30 1460 non-null   float64
 244  BsmtFinType1_31 1460 non-null   object 
 245  BsmtFinType2_31 1460 non-null   object 
 246  BsmtCond_31    1460 non-null   object 
 247  BsmtExposure_31 1460 non-null   object 
 248  BsmtExposure2_31 1460 non-null   object 
 249  BsmtFinPct_31 1460 non-null   float64
 250  BsmtUnfPct_31 1460 non-null   float64
 251  BsmtFinType1_32 1460 non-null   object 
 252  BsmtFinType2_32 1460 non-null   object 
 253  BsmtCond_32    1460 non-null   object 
 254  BsmtExposure_32 1460 non-null   object 
 255  BsmtExposure2_32 1460 non-null   object 
 256  BsmtFinPct_32 1460 non-null   float64
 257  BsmtUnfPct_32 1460 non-null   float64
 258  BsmtFinType1_33 1460 non-null   object 
 259  BsmtFinType2_33 1460 non-null   object 
 260  BsmtCond_33    1460 non-null   object 
 261  BsmtExposure_33 1460 non-null   object 
 262  BsmtExposure2_33 1460 non-null   object 
 263  BsmtFinPct_33 1460 non-null   float64
 264  BsmtUnfPct_33 1460 non-null   float64
 265  BsmtFinType1_34 1460 non-null   object 
 266  BsmtFinType2_34 1460 non-null   object 
 267  BsmtCond_34    1460 non-null   object 
 268  BsmtExposure_34 1460 non-null   object 
 269  BsmtExposure2_34 1460 non-null   object 
 270  BsmtFinPct_34 1460 non-null   float64
 271  BsmtUnfPct_34 1460 non-null   float64
 272  BsmtFinType1_35 1460 non-null   object 
 273  BsmtFinType2_35 1460 non-null   object 
 274  BsmtCond_35    1460 non-null   object 
 275  BsmtExposure_35 1460 non-null   object 
 276  BsmtExposure2_35 1460 non-null   object 
 277  BsmtFinPct_35 1460 non-null   float64
 278  BsmtUnfPct_35 1460 non-null   float64
 279  BsmtFinType1_36 1460 non-null   object 
 280  BsmtFinType2_36 1460 non-null   object 
 281  BsmtCond_36    1460 non-null   object 
 282  BsmtExposure_36 1460 non-null   object 
 283  BsmtExposure2_36 1460 non-null   object 
 284  BsmtFinPct_36 1460 non-null   float64
 285  BsmtUnfPct_36 1460 non-null   float64
 286  BsmtFinType1_37 1460 non-null   object 
 287  BsmtFinType2_37 1460 non-null   object 
 288  BsmtCond_37    1460 non-null   object 
 289  BsmtExposure_37 1460 non-null   object 
 290  BsmtExposure2_37 1460 non-null   object 
 291  BsmtFinPct_37 1460 non-null   float64
 292  BsmtUnfPct_37 1460 non-null   float64
 293  BsmtFinType1_38 1460 non-null   object 
 294  BsmtFinType2_38 1460 non-null   object 
 295  BsmtCond_38    1460 non-null   object 
 296  BsmtExposure_38 1460 non-null   object 
 297  BsmtExposure2_38 1460 non-null   object 
 298  BsmtFinPct_38 1460 non-null   float64
 299  BsmtUnfPct_38 1460 non-null   float64
 300  BsmtFinType1_39 1460 non-null   object 
 301  BsmtFinType2_39 1460 non-null   object 
 302  BsmtCond_39    1460 non-null   object 
 303  BsmtExposure_39 1460 non-null   object 
 304  BsmtExposure2_39 1460 non-null   object 
 305  BsmtFinPct_39 1460 non-null   float64
 306  BsmtUnfPct_39 1460 non-null   float64
 307  BsmtFinType1_40 1460 non-null   object 
 308  BsmtFinType2_40 1460 non-null   object 
 309  BsmtCond_40    1460 non-null   object 
 310  BsmtExposure_
```

```
In [122]: plt.figure(figsize=(15,10))
sns.heatmap(GD.corr())
plt.show()
```



```
In [123]: sns.pairplot(GD)
plt.show()
```



```
In [124]:
```

```
In [125]:
```

```
In [126]:
```

```

In [2]: import numpy as np
import pandas as pd
from pandas import DataFrame, Series
import matplotlib
import matplotlib.pyplot as plt
import statsmodels
import warnings
warnings.filterwarnings('ignore')
from sklearn import preprocessing
matplotlib.style.use('ggplot')

In [3]: from sklearn.datasets import load_boston

In [4]: boston=load_boston()

In [5]: print(boston.keys())
dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename'])

In [6]: print(boston.data.shape)
(506, 13)

In [7]: print(boston.feature_names)
['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO'
 'B' 'LSTAT']

In [8]: print(boston.DESCR)
...
.. boston_dataset:
Boston house prices dataset
-----
** Data Set Characteristics:**
 :Number of Instances: 506
 :Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target
 :Attribute Information (in order):
 - CRIM: per capita crime rate by town
 - ZN: proportion of residential land zoned for lots over 25,000 sq.ft.
 - INDUS: proportion of non-retail business acres per town
 - CHAS: Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
 - NOX: nitric oxide concentration (parts per 10 million)
 - RM: average number of rooms per dwelling
 - AGE: proportion of owner-occupied units built prior to 1940
 - DIS: weighted distances to five Boston employment centers
 - RAD: index of accessibility to radial highways
 - TAX: full-value property-tax rate per $10,000
 - PTRATIO: pupil-teacher ratio by town
 - B: 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
 - LSTAT: lower status of the population
 - MEDV: median value of owner-occupied homes in $1000's

:Missing Attribute Values: None
:Creator: Harrison, D. and Rubinfeld, D.L.
This is a copy of UCI ML housing dataset.
https://archive.ics.uci.edu/ml/machine-learning-databases/housing/
This data set was taken from the Statlib library which is maintained at Carnegie Mellon University.
The Boston house-price data of Harrison, J., D. et al. "Hedonic
Price and the Demand for Clean Air", J. Environ. Economics & Management,
vol. 1, pp. 81-102, 1978. Used in Belsley, Kuh & Welsh, "Regression Diagnostics
- Identifying Influential Data and Sources of Collinearity", 1980, pp. 244-245 of the latter.
The Boston house-price data has been used in many machine learning papers that address regression
problems.
.. topic:: References
- Belsley, Kuh, and Welsch, "Regression Diagnostics: Identifying Influential Data and Sources of Collinearity",
Wiley, 1980, 24-61.
- Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings of the Tenth International Conference on Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.

```

In [9]: bos=DataFrame(boston.data)

In [10]: bos.head()

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33

In [11]: bos.columns=boston.feature\_names

In [12]: bos

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33

506 rows × 13 columns

In [13]: bos['Price']=boston.target

In [14]: bos

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	Price
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

506 rows × 14 columns

In [15]: bos.describe()

In [16]: bos.info()

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	Price
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.
mean	3.613524	11.363636	11.16779	0.069170	0.554695	6.284634	68.57901	3.79504	9.549407	40.237154	18.455334	356.	24.0	
std	8.601545	23.322453	6.803533	0.253994	0.115878	0.02617	28.148861	2.105710	8.707259	168.537116	2.164546	9.1		
min	0.000000	0.000000	0.000000	0.000000	0.000000	5.000000	12.600000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	
25%	17.023000	5.812461	18.000000	0.000000	0.000000	5.000000	12.750000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	
50%	21.200000	18.000000	18.000000	0.000000	0.000000	5.000000	12.750000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	
75%	32.250000	22.000000	22.000000	0.000000	0.000000	5.000000	12.750000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	
max	50.000000	50.000000	50.000000	0.000000	0.000000	5.000000	12.750000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

In [17]: bos.info()
Name: Price, dtype: float64

In [18]: bos.describe()
Name: Price, dtype: float64

In [19]: plt.plot(bos['Price'])
plt.show()
sns.distplot(bos['Price'])
plt.show()
print('skewness={:.4f}'.format(bos['Price'].skew()))
print('kurtosis={:.4f}'.format(bos['Price'].kurt()))

skewness=1.08984028549072  
kurtosis=1.49519694165818

In [20]: bos.head(2)

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	Price
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

506 rows × 14 columns

In [21]: linear\_regress.describe()
Name: residuals, dtype: float64

In [22]: y=bos.iloc[:,13].values
y

In [23]: array([6.3200e-03, 1.8000e+01, 2.3100e+00, ..., 1.5300e+01, 3.9690e+02, 4.9800e+02, 8.0000e+02, 1.2000e+03, 1.6000e+03, 2.0000e+03, 2.4000e+03, 2.8000e+03, 3.2000e+03, 3.6000e+03, 4.0000e+03, 4.4000e+03, 4.8000e+03, 5.2000e+03, 5.6000e+03, 6.0000e+03, 6.4000e+03, 6.8000e+03, 7.2000e+03, 7.6000e+03, 8.0000e+03, 8.4000e+03, 8.8000e+03, 9.2000e+03, 9.6000e+03, 1.0000e+04, 1.0400e+04, 1.0800e+04, 1.1200e+04, 1.1600e+04, 1.2000e+04, 1.2400e+04, 1.2800e+04, 1.3200e+04, 1.3600e+04, 1.4000e+04, 1.4400e+04, 1.4800e+04, 1.5200e+04, 1.5600e+04, 1.6000e+04, 1.6400e+04, 1.6800e+04, 1.7200e+04, 1.7600e+04, 1.8000e+04, 1.8400e+04, 1.8800e+04, 1.9200e+04, 1.9600e+04, 2.0000e+04, 2.0400e+04, 2.0800e+04, 2.1200e+04, 2.1600e+04, 2.2000e+04, 2.2400e+04, 2.2800e+04, 2.3200e+04, 2.3600e+04, 2.4000e+04, 2.4400e+04, 2.4800e+04, 2.5200e+04, 2.5600e+04, 2.6000e+04, 2.6400e+04, 2.6800e+04, 2.7200e+04, 2.7600e+04, 2.8000e+04, 2.8400e+04, 2.8800e+04, 2.9200e+04, 2.9600e+04, 3.0000e+04, 3.0400e+04, 3.0800e+04, 3.1200e+04, 3.1600e+04, 3.2000e+04, 3.2400e+04, 3.2800e+04, 3.3200e+04, 3.3600e+04, 3.4000e+04, 3.4400e+04, 3.4800e+04, 3.5200e+04, 3.5600e+04, 3.6000e+04, 3.6400e+04, 3.6800e+04, 3.7200e+04, 3.7600e+04, 3.8000e+04, 3.8400e+04, 3.8800e+04, 3.9200e+04, 3.9600e+04, 4.0000e+04, 4.0400e+04, 4.0800e+04, 4.1200e+04, 4.1600e+04, 4.2000e+04, 4.2400e+04, 4.2800e+04, 4.3200e+04, 4.3600e+04, 4.4000e+04, 4.4400e+04, 4.4800e+04, 4.5200e+04, 4.5600e+04, 4.6000e+04, 4.6400e+04, 4.6800e+04, 4.7200e+04, 4.7600e+04, 4.8000e+04, 4.8400e+04, 4.8800e+04, 4.9200e+04, 4.9600e+04, 5.0000e+04, 5.0400e+04, 5.0800e+04, 5.1200e+04, 5.1600e+04, 5.2000e+04, 5.2400e+04, 5.2800e+04, 5.3200e+04, 5.3600e+04, 5.4000e+04, 5.4400e+04, 5.4800e+04, 5.5200e+04, 5.5600e+04, 5.6000e+04, 5.6400e+04, 5.6800e+04, 5.7200e+04, 5.7600e+04, 5.8000e+0



```
In [7]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [8]: data = pd.read_csv('Ecommerce_Customers.csv')

In [9]: data.head()

Out[9]:
```

	Email	Address	Avatar	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
0	mstephenson@fernandez.com	835 Frank Tunnel\nWrightmouth, MI 82180-9605	Violet	34.497268	12.655651	39.577668	4.082621	587.951054
1	hduke@hotmail.com	4547 Archer Common\nDiazchester, CA 06566-8576	DarkGreen	31.926272	11.109461	37.268959	2.664034	392.204933
2	pallen@yahoo.com	24645 Valerie Union Suite 582\nnCobbborough, D...	Bisque	33.000915	11.330278	37.110597	4.104543	487.547505
3	riverarebecca@gmail.com	1414 David Thoroughway\nnPort Jason, OH 22070-1220	SaddleBrown	34.305557	13.717514	36.721283	3.120179	581.852344
4	mstephens@davidson-herman.com	14023 Rodriguez Passage\nnPort Jacobville, PR 3...	MediumAquaMarine	33.330673	12.795189	37.536653	4.446308	599.406092

```
In [10]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
Email      500 non-null object
Address    500 non-null object
Avatar     500 non-null object
Avg. Session Length 500 non-null float64
Time on App   500 non-null float64
Time on Website 500 non-null float64
Length of Membership 500 non-null float64
Yearly Amount Spent 500 non-null float64
dtypes: float64(5), object(3)
memory usage: 25.5+ KB

In [11]: data.describe()
```

```
Out[11]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
count	500.000000	500.000000	500.000000	500.000000	500.000000
mean	33.053194	12.052488	37.060445	3.533462	499.314038
std	0.992563	0.994216	1.010489	0.999278	79.314782
min	29.532429	8.508152	33.913847	0.269901	256.670582
25%	32.341822	11.388153	36.349257	2.930450	445.038277
50%	33.082008	11.983231	37.069367	3.533975	498.887875
75%	33.711985	12.753850	37.716432	4.126502	549.313828
max	36.139662	15.126994	40.005182	6.922689	765.518462

```
Use seaborn to create a jointplot to compare the Time on Website and Yearly Amount Spent columns. Does the correlation make sense?
```

```
In [14]: plt.figure(figsize=(15,8),dpi=150)
sns.jointplot(x='Time on Website',y='Yearly Amount Spent',data=data)
plt.show()

<Figure size 2250x1200 with 0 Axes>
```

```
In [15]: data.corr()
```

```
Out[15]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
Avg. Session Length	1.000000	-0.027826	-0.034987	0.060247	0.355088
Time on App	-0.027826	1.000000	0.082388	0.029143	0.499328
Time on Website	-0.034987	0.082388	1.000000	-0.047582	-0.002641
Length of Membership	0.060247	0.029143	-0.047582	1.000000	0.809084
Yearly Amount Spent	0.355088	0.499328	-0.002641	0.809084	1.000000

```
In [16]: plt.figure(figsize=(15,8),dpi=150)
sns.jointplot(x='Time on App',y='Yearly Amount Spent',data=data)
plt.show()

<Figure size 2250x1200 with 0 Axes>
```

```
In [18]: plt.figure(figsize=(15,8),dpi=150)
sns.jointplot(x='Time on App',y='Length of Membership',data=data,kindle='hex')
plt.show()

<Figure size 2250x1200 with 0 Axes>
```

```
In [19]: plt.figure(figsize=(15,9),dpi=100)
sns.pairplot(data)
plt.show()

<Figure size 1500x900 with 0 Axes>
```

```
In [20]: sns.lmplot(x='Yearly Amount Spent',y='Length of Membership',data=data)
plt.show()

7
6
5
4
3
2
1
0
```

```
In [21]: data.columns
```

```
Out[21]: Index(['Email', 'Address', 'Avatar', 'Avg. Session Length', 'Time on App', 'Time on Website', 'Length of Membership', 'Yearly Amount Spent'], dtype='object')
```

```
In x used ['Avg. Session Length','Time on App','Time on Website','Length of Membership'] this columns of values
```

```
In [35]: x = data.drop(columns=['Yearly Amount Spent','Email','Address','Avatar'])

y = data['Yearly Amount Spent']

from sklearn.model_selection import train_test_split
```

```
In [38]: x_train,x_test,y_train,y_test = train_test_split(x,y,train_size=.80,random_state=0)
```

```
In [39]: x_train.shape,y_train.shape
```

```
Out[39]: ((400, 4), (400,))
```

```
In [40]: y_train.shape,y_test.shape
```

```
Out[40]: ((400,), (100,))
```

```
In [41]: from sklearn.linear_model import LinearRegression
```

```
In [42]: model = LinearRegression()
model.fit(x_train,y_train)

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [44]: y_pred=model.predict(x_test)
y_pred
```

```
Out[44]: array([438.46488066, 489.6618454, 370.06954186, 513.8590556, 495.69799869, 450.88726525, 458.01303756, 336.6655623, 528.5154919, 463.02782042, 559.42248461, 526.02431274, 552.20851747, 457.67586182, 450.15607779, 431.65097987, 563.58511617, 408.35440214, 599.61853986, 486.84826737, 716.65493547, 496.79327186, 411.49845257, 595.15395642, 551.84576857, 398.18338551, 421.30757507, 505.52726976, 546.87976787, 403.4873886, 558.43909636, 401.07293164, 503.60331067, 430.11282821, 569.6374307, 421.89472114, 487.26750654, 602.37423872, 616.44671945, 569.59491734, 449.73449152, 379.92615446, 558.4473283, 462.7200747, 502.16865444, 401.43460117, 513.84624174, 416.80395476, 576.471219, 478.50455211, 457.248031, 397.6712392, 564.8697802, 305.13130981, 63.075725, 476.3492249, 499.0418098, 433.98746406, 463.72067764, 516.91617781, 504.0395793, 356.62137493, 508.85162889, 416.95315655, 634.98001693, 455.88862211, 451.85573774, 767.06469393, 518.9651029, 457.34831524, 379.85834649, 375.07959328, 598.72889754, 515.75175435, 513.36496692, 429.46162916, 491.12451373, 491.58530432, 532.98278557, 753.43951807, 441.61268032, 507.23551171, 558.93278557, 523.33224278, 395.91913096, 516.9802983, 518.31262767, 383.74450882, 524.16819282, 500.68268787, 606.78782704, 408.80519096])
```

```
In [45]: from sklearn.metrics import r2_score
```

```
In [46]: print('R2 Score :',r2_score(y_test,y_pred))

R2 Score : 0.9861924261981548
```

```
In [47]: model.coef_
```

```
Out[47]: array([25.88815047, 38.87046474, 0.47066154, 61.78369022])
```

```
Create a scatterplot of the real test values versus the predicted values.
```

```
In [48]: sns.scatterplot(x=y_test,y=y_pred)
plt.show()

700
600
500
400
300
```

```
In [51]: # calculate these metrics by hand!
from sklearn import metrics

print('MAE:', metrics.mean_absolute_error(y_test, y_pred))
print('MSE:', metrics.mean_squared_error(y_test, y_pred))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))

MAE: 7.645674798915279
MSE: 92.88010304498495
RMSE: 9.63795118502812
```

```
Residuals
```

```
You should have gotten a very good model with a good fit. Let's quickly explore the residuals to make sure everything was okay with our data.
```

```
In [66]: sns.set_style('whitegrid')
err1=metrics.mean_absolute_error(y_test, y_pred)
err2=metrics.mean_squared_error(y_test, y_pred)
err3=np.sqrt(metrics.mean_squared_error(y_test, y_pred))

error =np.array([err1,err2,err3])

sns.distplot(data['Yearly Amount Spent'])
```

```
Out[67]: <matplotlib.axes._subplots.AxesSubplot at 0x146fb330>
```

```
In [92]: dataframe = pd.DataFrame({'Coefficient':model.coef_,index=['Avg. Session Length','Time on App','Time on Website','Length of Membership']})
```

```
Out[93]:
```

Coefficient
Avg. Session Length 25.888150
Time on App 38.870465
Time on Website 0.470662
Length of Membership 61.783690

```
Regression Coefficient.
```

```
Regression coefficients are estimates of the unknown population parameters and describe the relationship between a predictor variable and the response. In linear regression, coefficients are the values that multiply the predictor values. Suppose you have the following regression equation:  $y = 3X + 5$ . In this equation,  $+3$  is the coefficient,  $X$  is the predictor, and  $+5$  is the constant.
```

```
The sign of each coefficient indicates the direction of the relationship between a predictor variable and the response variable.
```

```
A positive sign indicates that as the predictor variable increases, the response variable also increases. A negative sign indicates that as the predictor variable increases, the response variable decreases. The coefficient value represents the mean change in the response given a one unit change in the predictor. For example, if a coefficient is  $+3$ , the mean response value increases by 3 for every one unit change in the predictor.
```

```
Do you think the company should focus more on their mobile app or on their website?
```

```
On their app
```

```
In [21]: data.columns
```

```
Out[21]: Index(['Email', 'Address', 'Avatar', 'Avg. Session Length', 'Time on App', 'Time on Website', 'Length of Membership', 'Yearly Amount Spent'], dtype='object')
```

```
In x used ['Avg. Session Length','Time on App','Time on Website','Length of Membership'] this columns of values
```

```
In [35]: x = data.drop(columns=['Yearly Amount Spent','Email','Address','Avatar'])

y = data['Yearly Amount Spent']

from sklearn.model_selection import train_test_split
```

```
In [38]: x_train,x_test,y_train,y_test = train_test_split(x,y,train_size=.80,random_state=0)
```

```
In [39]: x_train.shape,y_train.shape
```

```
Out[39]: ((400, 4), (400,))
```

```
In [40]: y_train.shape,y_test.shape
```

```
Out[40]: ((400,), (100,))
```

```
In [41]: from sklearn.linear_model import LinearRegression
```

```
In [42]: model = LinearRegression()
model.fit(x_train,y_train)

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [44]: y_pred=model.predict(x_test)
y_pred
```

```
Out[44]: array([438.46488066, 489.6618454, 370.06954186, 513.8590556, 495.69799869, 450.88726525, 458.01303756, 336.6655623, 528.5154919, 463.02782042, 559.42248461, 526.02431274, 552.20851747, 457.67586182, 450.15607779, 431.65097987, 563.58511617, 408.35440214, 599.61853986, 486.84826737, 716.65493547, 496.79327186, 411.49845257, 595.15395642, 551.84576857, 398.18338551, 421.30757507, 505.52726976, 546.8697802, 305.13130981, 63.075725, 476.3492249, 499.0418098, 433.98746406, 463.72067764, 516.91617781, 504.0395793, 356.62137493, 508.85162889, 416.95315655, 634.98001693, 455.88862211, 451.85573774, 767.06469393, 518.9651029, 457.34831524, 379.85834649, 375.07959328, 598.72889754, 515.75175435, 513.36496692, 429.46162916, 491.12451373, 491.58530432, 532.98278557, 753.43951807, 441.61268032, 507.23551171, 558.93278557, 523.33224278, 395.91913096, 516.9802983, 518.31262767, 383.74450882, 524.16819282, 500.68268787, 606.78782704, 408.80519096])
```

```
In [45]: from sklearn.metrics import r2_score
```

```
In [46]: print('R2 Score :',r2_score(y_test,y_pred))

R2 Score : 0.9861924261981548
```

```
In [47]: model.coef_
```

```
Out[47]: array([25.88815047, 38.87046474, 0.47066154, 61.78369022])
```

```
Create a scatterplot of the real test values versus the predicted values.
```

```
In [48]: sns.scatterplot(x=y_test,y=y_pred)
plt.show()

700
600
500
400
300
```

```
In [51]: # calculate these metrics by hand!
from sklearn import metrics

print('MAE:', metrics.mean_absolute_error(y_test, y_pred))
print('MSE:', metrics.mean_squared_error(y_test, y_pred))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))

MAE: 7.645674798915279
MSE: 92.88010304498495
RMSE: 9.63795118502812
```

```
Residuals
```

```
You should have gotten a very good model with a good fit. Let's quickly explore the residuals to make sure everything was okay with our data.
```

```
In [66]: sns.set_style('whitegrid')
err1=metrics.mean_absolute_error(y_test, y_pred)
err2=metrics.mean_squared_error(y_test, y_pred)
err3=np.sqrt(metrics.mean_squared_error(y_test, y_pred))

error =np.array([err1,err2,err3])

sns.distplot(data['Yearly Amount Spent'])
```

```
Out[67]: <matplotlib.axes._subplots.AxesSubplot at 0x146fb330>
```

```
In [92]: dataframe = pd.DataFrame({'Coefficient':model.coef_,index=['Avg. Session Length','Time on App','Time on Website']})
```

```
Out[93]:
```

Coefficient
Avg. Session Length 25.888150
Time on App 38.870465
Time on Website 0.470662

```
Regression Coefficient.
```

```
Regression coefficients are estimates of the unknown population parameters and describe the relationship between a predictor variable and the response. In linear regression, coefficients are the values that multiply the predictor values. Suppose you have the following regression equation:  $y = 3X + 5$ . In this equation,  $+3$  is the coefficient,  $X$  is the predictor, and  $+5$  is the constant.
```

```
The sign of each coefficient indicates the direction of the relationship between a predictor variable and the response variable.
```

```
A positive sign indicates that as the predictor variable increases, the response variable also increases. A negative sign indicates that as the predictor variable increases, the response variable decreases. The coefficient value represents the mean change in the response given a one unit change in the predictor. For example, if a coefficient is  $+3$ , the mean response value increases by 3 for every one unit change in the predictor.
```

```
Do you think the company should focus more on their mobile app or on their website?
```

```
On their app
```



```
In [157]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
sns.set(style="darkgrid")

In [158]: from sklearn.datasets import load_boston

In [159]: boston = load_boston()

In [160]: print(boston.keys())
dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename'])

In [161]: print(boston.data.shape)
(506, 13)

In [162]: print(boston.feature_name)
['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO'
 'B' 'LSTAT']

In [163]: print(boston.DESCR)

Boston house-prices dataset
as always passes asssessments

**Data Sets Characteristics**
:Number of Instances: 506
:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target
:t:Median value of owner-occupied homes in $1000's

:Attribute Information (in order):
 - CRIM : per capita crime rate by town
 - ZN : proportion of residential land zoned for lots over 25,000 sq.ft.
 - INDUS : proportion of non-retail business acres per town
 - CHAS : Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
 - NOX : nitric oxide concentration (parts per 10,000)
 - RM : average number of rooms per dwelling
 - AGE : proportion of owner-occupied units built prior to 1940
 - DIS : weighted distances to five Boston employment centres
 - RAD : index of accessibility to radial highways
 - TAX : full-value property-tax rate per $10,000
 - PTRATIO : pupil-teacher ratio by town
 - B : lower status of the population
 - LSTAT : % lower status of the population
 - MEDV : Median value of owner-occupied homes in $1000's

:Missing Attribute Values: None
:Creator: Harrison, D. and Rubinfeld, D.L.
This is a copy of UCI ML housing dataset.
https://archive.ics.uci.edu/ml/machine-learning-databases/housing/

This dataset was taken from the Statlib library which is maintained at Carnegie Mellon University.
The Boston house-prices data of Harrison, D. and Rubinfeld, D.L., 'Hedonic prices and the demand for clean air', J. Environ. Economics & Management, vol.5, 81-102, 1978. N.B. Used in Nelson, Kuh & Welch, 'Regression diagnostics ...', Wiley, 1980. Various transformations are used in the table on pages 244-261 of the latter.

The Boston house-prices data has been used in many machine learning papers that address regression problems.

.. topic:: References
 - Belshay, Kuh & Welch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980, 244-261.
 - Quinlan,R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings of the Tenth International Conference on Machine Learning. Massachusetts, Amherst. Morgan Kaufmann.

In [164]: data = pd.DataFrame(boston.data,columns=boston.feature_names)

In [165]: data.head()
Out[165]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT
0 0.0652 18.0 2.31 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
1 0.02731 0.0 7.07 0.0 0.469 6.421 78.9 4.9671 2.0 242.0 17.8 396.90 9.14
2 0.02729 0.0 7.07 0.0 0.469 7.185 61.1 4.9671 2.0 242.0 17.8 392.83 4.03
3 0.03237 0.0 2.18 0.0 0.458 6.998 458.60622 3.0 222.0 18.7 394.63 2.94
4 0.06905 0.0 2.18 0.0 0.458 7.147 542.60622 3.0 222.0 18.7 396.90 5.33

In [166]: data.tail()
Out[166]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT
501 0.06263 0.0 11.93 0.0 0.573 6.593 69.1 2.4786 1.0 273.0 21.0 391.99 9.57
502 0.04527 0.0 11.93 0.0 0.573 6.120 76.7 2.2875 1.0 273.0 21.0 396.90 9.08
503 0.06076 0.0 11.93 0.0 0.573 6.976 91.0 2.1675 1.0 273.0 21.0 396.90 5.64
504 0.10595 0.0 11.93 0.0 0.573 6.794 89.3 2.3889 1.0 273.0 21.0 393.45 6.48
505 0.04741 0.0 11.93 0.0 0.573 6.630 80.8 2.5050 1.0 273.0 21.0 396.90 7.88

In [167]: data[['Price']] = boston.target

In [168]: data.head()
Out[168]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT Price
0 0.0652 18.0 2.31 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98 240
1 0.02731 0.0 7.07 0.0 0.469 6.421 78.9 4.9671 2.0 242.0 17.8 396.90 9.14 216
2 0.02729 0.0 7.07 0.0 0.469 7.185 61.1 4.9671 2.0 242.0 17.8 392.83 4.03 347
3 0.03237 0.0 2.18 0.0 0.458 6.998 458.60622 3.0 222.0 18.7 394.63 2.94 334
4 0.06905 0.0 2.18 0.0 0.458 7.147 542.60622 3.0 222.0 18.7 396.90 5.33 362

Inferential statics.
EDA
Outlier treatment apply on target.
Apply Machine Algorithm i.e. Linear regression here.

In [169]: data.describe()
Out[169]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT
count 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000
mean 3.6495 23.322436 11.363453 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
std 8.601545 20.000000 0.53828 0.042597 0.427321 0.693564 0.693564 0.648536 0.648536 0.648536 0.648536 0.648536
min 0.06520 0.000000 0.46900 0.000000 0.385000 3.56100 2.90000 1.129600 1.000000 1.000000 1.000000 1.000000
25% 0.25650 0.000000 0.59000 0.000000 0.58000 5.08000 5.08000 0.429600 0.429600 0.429600 0.429600 0.429600
50% 0.3677083 12.500000 18.100000 0.000000 0.62400 6.623500 6.623500 0.970500 0.970500 0.970500 0.970500 0.970500
75% 0.4777083 20.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000
max 88.97602 100.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000

In [170]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 14 columns):
CRIM 320 non-null float64
ZN 506 non-null float64
INDUS 506 non-null float64
CHAS 506 non-null float64
NOX 506 non-null float64
RM 506 non-null float64
AGE 506 non-null float64
DIS 506 non-null float64
RAD 506 non-null float64
TAX 506 non-null float64
PTRATIO 506 non-null float64
B 506 non-null float64
LSTAT 506 non-null float64
Price 506 non-null float64
dtype: float64 (nbytes: 55.4 KB)

In [171]: data.sum()
Out[171]: CRIM 0
ZN 0
INDUS 0
CHAS 0
NOX 0
RM 0
AGE 0
DIS 0
RAD 0
TAX 0
PTRATIO 0
B 0
LSTAT 0
Price 0

In [172]: data.duplicated().sum()
Out[172]: 0

EDA
In [173]: plt.figure(figsize=(15,9),dpi=200)
sns.pairplot(data)
plt.show()

Outlier treatment apply on target.
Apply Machine Algorithm i.e. Linear regression here.

In [174]: data.describe()
Out[174]: count 506.000000
mean 3.6495 23.322436 11.363453 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
std 8.601545 20.000000 0.53828 0.042597 0.427321 0.693564 0.693564 0.648536 0.648536 0.648536
min 0.06520 0.000000 0.46900 0.000000 0.385000 3.56100 2.90000 1.129600 1.000000 1.000000 1.000000 1.000000
25% 0.25650 0.000000 0.59000 0.000000 0.58000 5.08000 5.08000 0.429600 0.429600 0.429600 0.429600 0.429600
50% 0.3677083 12.500000 18.100000 0.000000 0.62400 6.623500 6.623500 0.970500 0.970500 0.970500 0.970500 0.970500
75% 0.4777083 20.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000
max 88.97602 100.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000

In [175]: data['Price'].describe()
Out[175]: count 506.000000
mean 32.52806
std 9.197104
min 0.000000
25% 8.125000
50% 21.200000
75% 35.000000
max 50.000000
Name: Price, dtype: float64

In [176]: np.abs(zscore(data['Price']))
Out[176]: array([ 1.519191,  1.519191,  1.519191, ...,  1.519191,  1.519191])

In [177]: data['z']=zscore(data['Price'])
Out[177]: data.head()
Out[177]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT z
0 0.0652 18.0 2.31 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98 240
1 0.02731 0.0 7.07 0.0 0.469 6.421 78.9 4.9671 2.0 242.0 17.8 396.90 9.14 216
2 0.02729 0.0 7.07 0.0 0.469 7.185 61.1 4.9671 2.0 242.0 17.8 392.83 4.03 347
3 0.03237 0.0 2.18 0.0 0.458 6.998 458.60622 3.0 222.0 18.7 394.63 2.94 334
4 0.06905 0.0 2.18 0.0 0.458 7.147 542.60622 3.0 222.0 18.7 396.90 5.33 362

Inferential statics.
EDA
Out[178]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT
count 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000
mean 3.6495 23.322436 11.363453 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
std 8.601545 20.000000 0.53828 0.042597 0.427321 0.693564 0.693564 0.648536 0.648536 0.648536 0.648536 0.648536
min 0.06520 0.000000 0.46900 0.000000 0.385000 3.56100 2.90000 1.129600 1.000000 1.000000 1.000000 1.000000
25% 0.25650 0.000000 0.59000 0.000000 0.58000 5.08000 5.08000 0.429600 0.429600 0.429600 0.429600 0.429600
50% 0.3677083 12.500000 18.100000 0.000000 0.62400 6.623500 6.623500 0.970500 0.970500 0.970500 0.970500 0.970500
75% 0.4777083 20.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000
max 88.97602 100.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000

In [179]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 14 columns):
CRIM 320 non-null float64
ZN 506 non-null float64
INDUS 506 non-null float64
CHAS 506 non-null float64
NOX 506 non-null float64
RM 506 non-null float64
AGE 506 non-null float64
DIS 506 non-null float64
RAD 506 non-null float64
TAX 506 non-null float64
PTRATIO 506 non-null float64
B 506 non-null float64
LSTAT 506 non-null float64
Price 506 non-null float64
dtype: float64 (nbytes: 55.4 KB)

In [180]: data.sum()
Out[180]: CRIM 0
ZN 0
INDUS 0
CHAS 0
NOX 0
RM 0
AGE 0
DIS 0
RAD 0
TAX 0
PTRATIO 0
B 0
LSTAT 0
Price 0

In [181]: data.duplicated().sum()
Out[181]: 0

EDA
In [182]: plt.figure(figsize=(15,9),dpi=200)
sns.pairplot(data)
plt.show()

Outlier treatment apply on target.
Apply Machine Algorithm i.e. Linear regression here.

In [183]: data.describe()
Out[183]: count 506.000000
mean 3.6495 23.322436 11.363453 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
std 8.601545 20.000000 0.53828 0.042597 0.427321 0.693564 0.693564 0.648536 0.648536 0.648536 0.648536 0.648536
min 0.06520 0.000000 0.46900 0.000000 0.385000 3.56100 2.90000 1.129600 1.000000 1.000000 1.000000 1.000000
25% 0.25650 0.000000 0.59000 0.000000 0.58000 5.08000 5.08000 0.429600 0.429600 0.429600 0.429600 0.429600
50% 0.3677083 12.500000 18.100000 0.000000 0.62400 6.623500 6.623500 0.970500 0.970500 0.970500 0.970500 0.970500
75% 0.4777083 20.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000
max 88.97602 100.000000 27.740000 1.000000 0.871000 8.780000 100.000000 12.126500 24.000000 71.000000 22.000000 396.000000

In [184]: data['Price'].describe()
Out[184]: count 506.000000
mean 32.52806
std 9.197104
min 0.000000
25% 8.125000
50% 21.200000
75% 35.000000
max 50.000000
Name: Price, dtype: float64

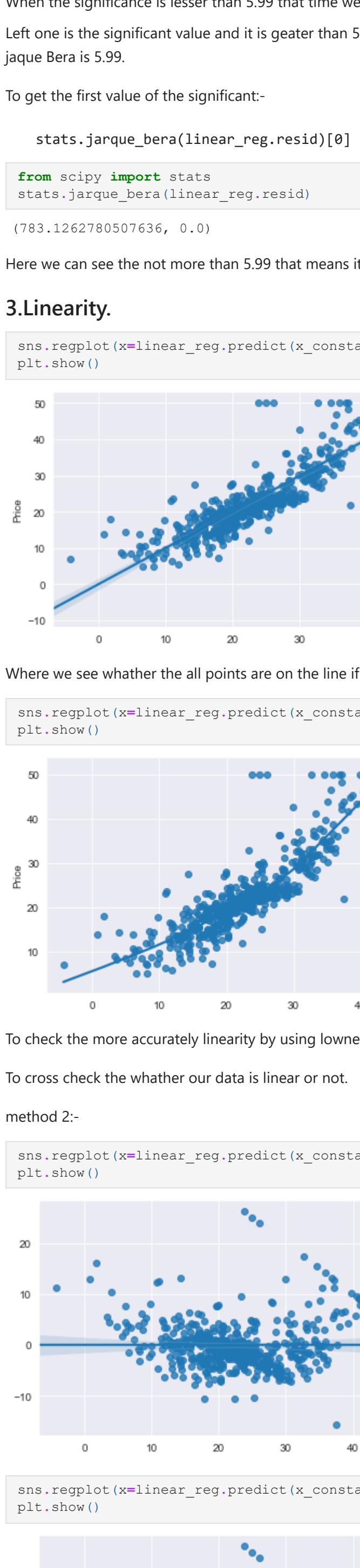
In [185]: np.abs(zscore(data['Price']))
Out[185]: array([ 1.519191,  1.519191,  1.519191, ...,  1.519191,  1.519191])

In [186]: data['z']=zscore(data['Price'])
Out[186]: data.head()
Out[186]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT z
0 0.0652 18.0 2.31 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98 240
1 0.02731 0.0 7.07 0.0 0.469 6.421 78.9 4.9671 2.0 242.0 17.8 396.90 9.14 216
2 0.02729 0.0 7.07 0.0 0.469 7.185 61.1 4.9671 2.0 242.0 17.8 392.83 4.03 347
3 0.03237 0.0 2.18 0.0 0.458 6.998 458.60622 3.0 222.0 18.7 394.63 2.94 334
4 0.06905 0.0 2.18 0.0 0.458 7.147 542.60622 3.0 222.0 18.7 396.90 5.33 362

Inferential statics.
EDA
Out[187]: CRIM ZN INDUS CHAS NOX RM AGE DIS RAD TAX PTRATIO B LSTAT
count 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000 506.000000
mean 3.6495 23.322436 11.363453 0.0 0.538 6.575 652.40900 1.0 296.0 15.3 396.90 4.98
std 8.601545 20.000000 0.53828 0.04259
```

## 2nd Assumptions: Multivariate Normality.

```
In [248]: sns.distplot(linear_reg.resid)
```



From the above we can see the graph is normally distributed.

### test2 :- Jarque Bera Test.

This test is for goodness of fit test of whether the sample data have skewness and kurtosis matching to the normally distributed or not.

Null Hypothesis:- Error should be normally distributed

Alternative hypothesis:- Error terms are not normally distributed.

When the significance is greater than 5.99 that time we can say error terms are normally distributed.

When the significance is lesser than 5.99 that time we can say the error terms are normally distributed.

Let the significance value and it is greater than 5.99(significant value) that means the it is normally distributed. significant value for jarque Bera is 5.99.

To get the first value of the significant:-

```
stats.jarque_bera(linear_reg.resid)[0]
```

```
In [249]: from scipy import stats
```

```
stats.jarque_bera(linear_reg.resid)
```

```
Out[249]: (783.126278057636, 0.0)
```

Here we can see the not more than 5.99 that means it is not normally distributed.

### 3. Linearity.

```
In [250]: sns.regplot(x=linear_reg.predict(x_constant), y=y)
```



Where we see whether all the points are on the line if yes then data is linear.

```
In [251]: sns.regplot(x=linear_reg.predict(x_constant), y=y, lowess=True)
```



### method 2:- RainBow Test

Null Hypothesis is the belongs to linear hypothesis. Alternate hypothesis is belongs to non-linear Hypothesis.

```
In [252]: sm.stats.diagnostic.linear_rainbow(res=linear_reg)
```

```
Out[252]: (0.5920850763961426, 0.99998852162534)
```

in the return we having the first value is f-test and second is the p-test. here we can see that the pvalue is greater than the 0.05 that means the model is linear.

here we can see that model is linear model.

### Method 2:-

if residual mean is close to the zero that means the the model is linear.

```
In [253]: np.mean(linear_reg.resid)
```

```
Out[253]: 1.227301089040173e-13
```

Since the residual mean is close to zero then it is linear.

### Homocidaicity

If the variance of the residual are symmetrically distributed across the residual line then data is said to homesidaicity.

If the variance is unequal residual across the residual line then data is said to heterosidaicity.

```
In [254]: sns.regplot(x=linear_reg.predict(x_constant), y=linear_reg.resid)
```



method 3:- GoldField Test or Brush Wagon test.

Homocidaicity:- When the variance is constant across the line belongs to null Hypothesis.

Heterocidaicity:- When the variance is not constant across line belongs to Alternate Hypothesis.

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

```
In [255]: import statsmodels.stats.api as sms
```

```
In [256]: sms.het_goldfeldquandt(linear_reg.resid, linear_reg.model.exog)
```

```
Out[256]: (2.4311473253497458, 7.302071915608868e-12, 'increasing')
```

```
In [257]: linear_reg.model.exog.shape
```

```
Out[257]: (506, 14)
```

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

### Multi-Colinearity

It found in adversely in our analysis. There should be No and low Multi-colinearity.

Correlation Causation:- Formulative relation between the two entity.

Causation:- Deterministic relationship between the data like temp in Oc and F(Farenite).

VIF - Variance Influential Factor. (to check the multi-colinearity)

VIF is the quotient of variance in model of multiple term by the variance with one term. It tells about the multi-colinearity that is ratio.

```
In [264]: from statsmodels.stats.outliers_influence import variance_inflation_factor
```

```
In [265]: vif = [variance_inflation_factor(x_constant.values, i) for i in range(x_constant.shape[1])]
```

```
In [266]: vif
```

```
Out[266]: [585.2652379423121, 1.227301089040173e-13, 2.2987581787494418, 3.991596183460297, 1.0739361275537886, 1.313961275537886, 1.933744357832574, 3.1100825128153372, 3.9908251281533725, 7.484493537474727, 9.00855394759707, 1.9908251281533725, 1.313961275537886, 2.9414910789319366]
```

```
In [275]: datafram=pd.DataFrame(vif,columns=["vif"],index=x_constant.columns)
```

```
Out[275]: vif
```



method 2:- GoldField Test or Brush Wagon test.

Homocidaicity:- When the variance is constant across the line belongs to null Hypothesis.

Heterocidaicity:- When the variance is not constant across line belongs to Alternate Hypothesis.

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

```
In [276]: import statsmodels.stats.api as sms
```

```
In [277]: sms.het_goldfeldquandt(linear_reg.resid, linear_reg.model.exog)
```

```
Out[277]: (2.4311473253497458, 7.302071915608868e-12, 'increasing')
```

```
In [278]: linear_reg.model.exog.shape
```

```
Out[278]: (506, 14)
```

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

### method 3:- RainBow Test

Null Hypothesis is the belongs to linear hypothesis. Alternate hypothesis is belongs to non-linear Hypothesis.

```
In [279]: sm.stats.diagnostic.linear_rainbow(res=linear_reg)
```

```
Out[279]: (0.5920850763961426, 0.99998852162534)
```

in the return we having the first value is f-test and second is the p-test. here we can see that the pvalue is greater than the 0.05 that means the model is linear.

here we can see that model is linear model.

### Method 2:-

if residual mean is close to the zero that means the the model is linear.

```
In [280]: np.mean(linear_reg.resid)
```

```
Out[280]: 1.227301089040173e-13
```



method 3:- GoldField Test or Brush Wagon test.

Homocidaicity:- When the variance is constant across the line belongs to null Hypothesis.

Heterocidaicity:- When the variance is not constant across line belongs to Alternate Hypothesis.

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

```
In [281]: import statsmodels.stats.api as sms
```

```
In [282]: sms.het_goldfeldquandt(linear_reg.resid, linear_reg.model.exog)
```

```
Out[282]: (2.4311473253497458, 7.302071915608868e-12, 'increasing')
```

```
In [283]: linear_reg.model.exog.shape
```

```
Out[283]: (506, 14)
```

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

### method 3:- RainBow Test

Null Hypothesis is the belongs to linear hypothesis. Alternate hypothesis is belongs to non-linear Hypothesis.

```
In [284]: sm.stats.diagnostic.linear_rainbow(res=linear_reg)
```

```
Out[284]: (0.5920850763961426, 0.99998852162534)
```

in the return we having the first value is f-test and second is the p-test. here we can see that the pvalue is greater than the 0.05 that means the model is linear.

here we can see that model is linear model.

### Method 2:-

if residual mean is close to the zero that means the the model is linear.

```
In [285]: np.mean(linear_reg.resid)
```

```
Out[285]: 1.227301089040173e-13
```



method 3:- GoldField Test or Brush Wagon test.

Homocidaicity:- When the variance is constant across the line belongs to null Hypothesis.

Heterocidaicity:- When the variance is not constant across line belongs to Alternate Hypothesis.

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

```
In [286]: import statsmodels.stats.api as sms
```

```
In [287]: sms.het_goldfeldquandt(linear_reg.resid, linear_reg.model.exog)
```

```
Out[287]: (2.4311473253497458, 7.302071915608868e-12, 'increasing')
```

```
In [288]: linear_reg.model.exog.shape
```

```
Out[288]: (506, 14)
```

It return the t-value,p\_value,order

Here we can see that pvalue is lesser than the 0.005 that means it rejected the null hypothesis and accepted the alternate hypothesis.it return the heterosidaicity.

### method 3:- RainBow Test

Null Hypothesis is the belongs to linear hypothesis. Alternate hypothesis is belongs to non-linear Hypothesis.

```
In [289]: sm.stats.diagnostic.linear_rainbow(res=linear_reg)
```

```
Out[289]: (0.5920850763961426, 0.99998852162534)
```

in the return we having the first value is f-test and second is the p-test. here we can see that the pvalue is greater than the 0.05 that means the model is linear.

here we can see that model is linear model.

### Method 2:-

if residual mean is close to the zero that means the the model is linear.

```
In [290]: np.mean(linear_reg.resid)
```

```
Out[290]: 1.227301089040173e-13
```



method 3:- GoldField Test or Brush Wagon test.