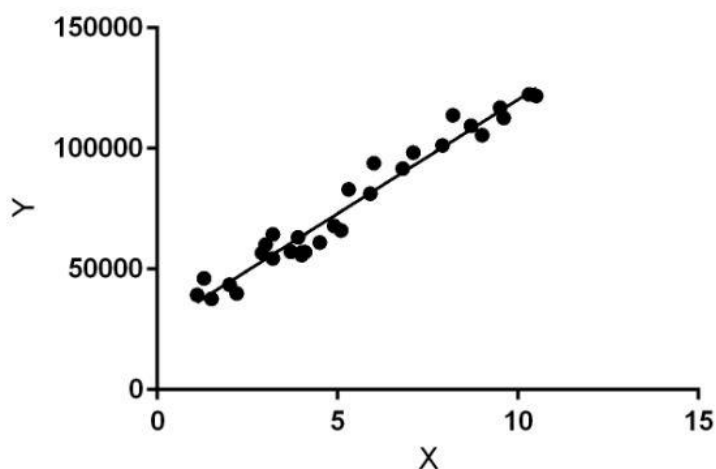| |
|---|
| Experiment No. 1 |
| Analyze the Boston Housing dataset and apply appropriate Regression Technique |
| Date of Performance: |
| Date of Submission: |

**Aim:**

Analyze the Boston Housing dataset and apply appropriate Regression Technique.

## Objective:

Ablility to perform various feature engineering tasks, apply linear regression on the given dataset and minimise the error.

## Theory:

Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables they are considering, and the number of independent variables getting used.

Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y(output). Hence, the name is Linear Regression.

In the figure above, X (input) is the work experience and Y (output) is the salary of a person. The regression line is the best fit line for our model.

## Dataset:

The Boston Housing Dataset

The Boston Housing Dataset is a derived from information collected by the U.S. Census Service concerning housing in the area of Boston MA. The following describes the dataset

## columns:

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 oxides 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 centres

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

## Code:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
BostonTrain = pd.read_csv("/content/archive (3).zip")
BostonTrain.head()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | MEDV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1 | 296 | 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 | 242 | 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 | 242 | 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 | 222 | 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 | 222 | 18.7 | 396.90 | NaN | 36.2 |

```python
BostonTrain.info()
BostonTrain.describe()
```

```
0   CRIM     486 non-null   float64
1   ZN       486 non-null   float64
2   INDUS    486 non-null   float64
3   CHAS     486 non-null   float64
4   NOX      506 non-null   float64
5   RM       506 non-null   float64
6   AGE      486 non-null   float64
7   DIS      506 non-null   float64
8   RAD      506 non-null   int64
9   TAX      506 non-null   int64
10  PTRATIO  506 non-null   float64
11  B        506 non-null   float64
12  LSTAT    486 non-null   float64
13  MEDV     506 non-null   float64
dtypes: float64(12), int64(2)
memory usage: 55.5 KB
```

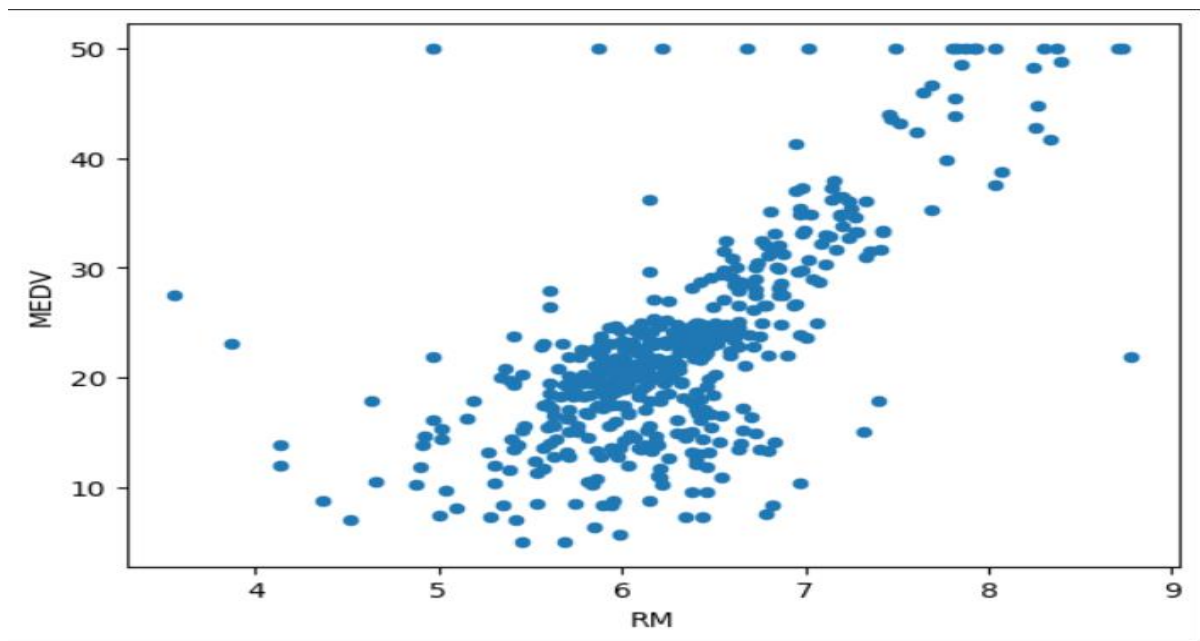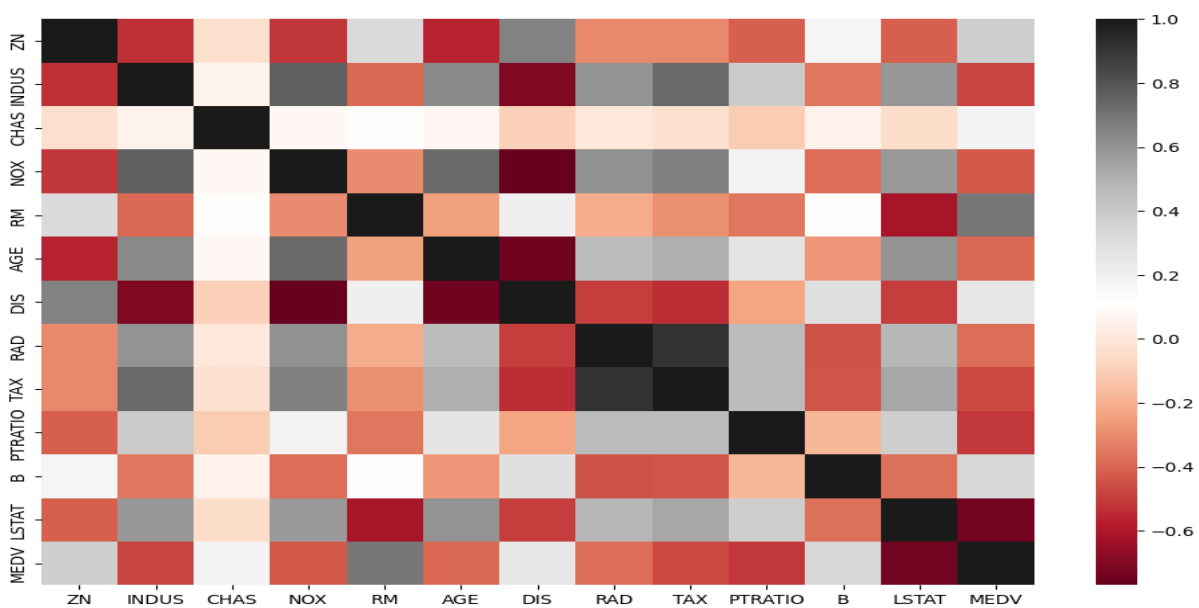| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | MEDV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 486.000000 | 486.000000 | 486.000000 | 486.000000 | 506.000000 | 506.000000 | 486.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 486.000000 | 506.000000 |
| mean | 3.611874 | 11.211934 | 11.083992 | 0.069959 | 0.554695 | 6.284634 | 68.518519 | 3.795043 | 9.549407 | 408.237154 | 18.455534 | 356.674032 | 12.715432 | 22.532806 |
| std | 8.720192 | 23.388876 | 6.835896 | 0.255340 | 0.115878 | 0.702617 | 27.999513 | 2.105710 | 8.707259 | 168.537116 | 2.164946 | 91.294864 | 7.155871 | 9.197104 |
| min | 0.006320 | 0.000000 | 0.460000 | 0.000000 | 0.385000 | 3.561000 | 2.900000 | 1.129600 | 1.000000 | 187.000000 | 12.600000 | 0.320000 | 1.730000 | 5.000000 |
| 25% | 0.081900 | 0.000000 | 5.190000 | 0.000000 | 0.449000 | 5.885500 | 45.175000 | 2.100175 | 4.000000 | 279.000000 | 17.400000 | 375.377500 | 7.125000 | 17.025000 |
| 50% | 0.253715 | 0.000000 | 9.690000 | 0.000000 | 0.538000 | 6.208500 | 76.800000 | 3.207450 | 5.000000 | 330.000000 | 19.050000 | 391.440000 | 11.430000 | 21.200000 |
| 75% | 3.560263 | 12.500000 | 18.100000 | 0.000000 | 0.624000 | 6.623500 | 93.975000 | 5.188425 | 24.000000 | 666.000000 | 20.200000 | 396.225000 | 16.955000 | 25.000000 |
| max | 88.976200 | 100.000000 | 27.740000 | 1.000000 | 0.871000 | 8.780000 | 100.000000 | 12.126500 | 24.000000 | 711.000000 | 22.000000 | 396.900000 | 37.970000 | 50.000000 |

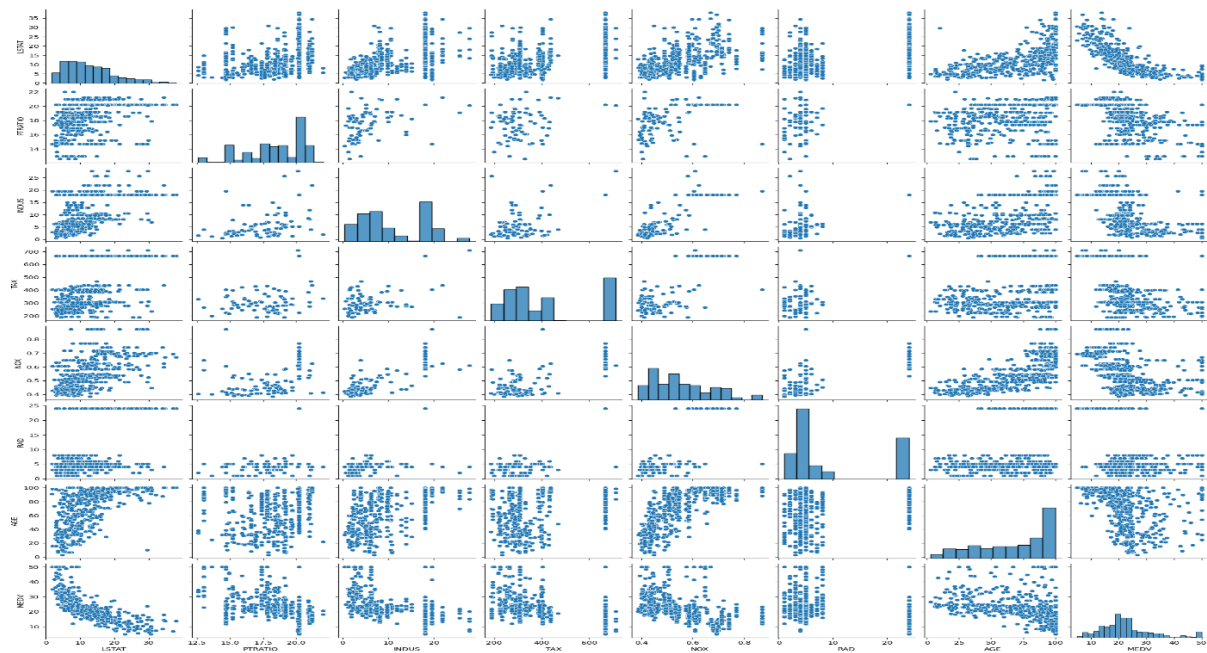BostonTrain.drop('CRIM', axis = 1, inplace=True)

BostonTrain.plot.scatter('RM', 'MEDV')



plt.subplots(figsize=(12,8))

sns.heatmap(BostonTrain.corr(), cmap = 'RdGy')



sns.pairplot(BostonTrain, vars = ['LSTAT', 'PTRATIO', 'INDUS', 'TAX', 'NOX', 'RAD', 'AGE', 'MEDV'])

## Conclusion:

What are features have been chosen to develop the model? Justify the features chosen to estimate the price of a house.

1. CRIM - per capita crime rate by town
2. CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise)
3. NOX - nitric oxides concentration (parts per 10 million)
4. RM - average number of rooms per dwelling
5. DIS - weighted distances to five Boston employment centres
6. RAD - index of accessibility to radial highways
7. TAX - full-value property-tax rate per $10,000
8. LSTAT - % lower status of the population

Comment on the Mean Squared Error calculated

1. Calculate Mean Square Error:0.04(+/- 0.04)
2. The Mean Square Error measures how close a regression line is to a set of data points
3. Lesser the Mean Square Error refers to Smaller is the error and Better the estimator