# LAB REPORT 6
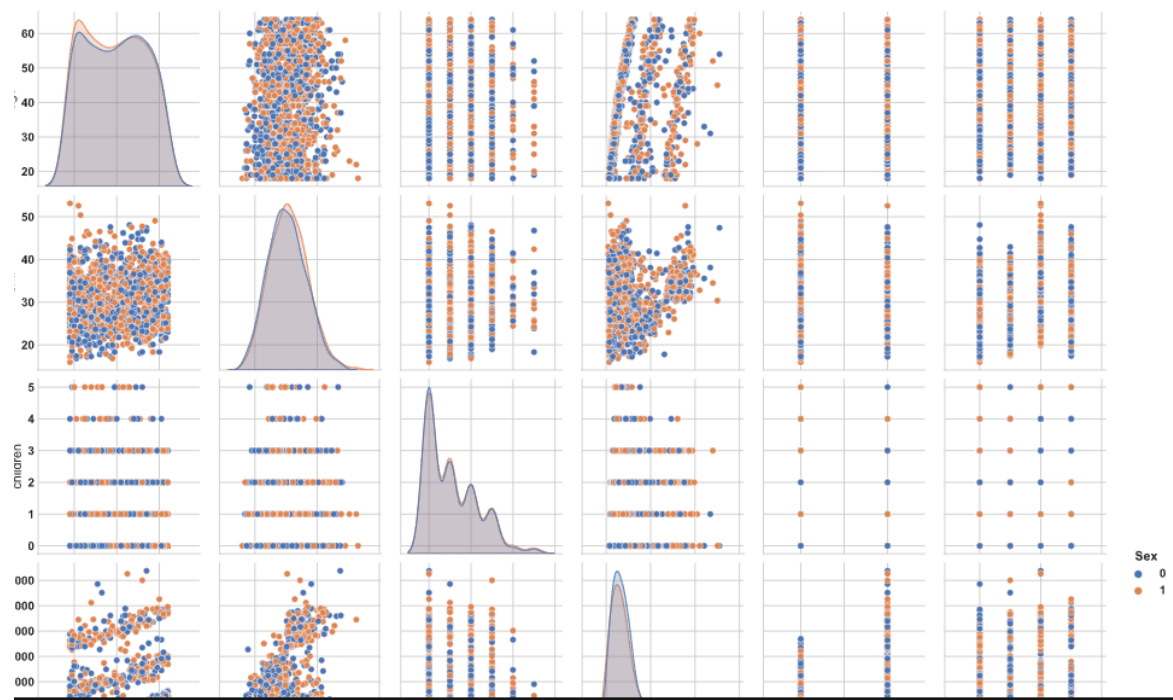
**Mayank Raj (B19CSE053)**

B.Tech CSE

## 1) Exploratory data analysis

```
    df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 10 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       1338 non-null   int64
 1   sex       1338 non-null   object
 2   bmi       1338 non-null   float64
 3   children  1338 non-null   int64
 4   smoker    1338 non-null   object
 5   region    1338 non-null   object
 6   charges   1338 non-null   float64
 7   Sex       1338 non-null   int32
 8   Smoker    1338 non-null   int32
 9   Region    1338 non-null   int32
dtypes: float64(2), int32(3), int64(2), object(3)
memory usage: 89.0+ KB
```
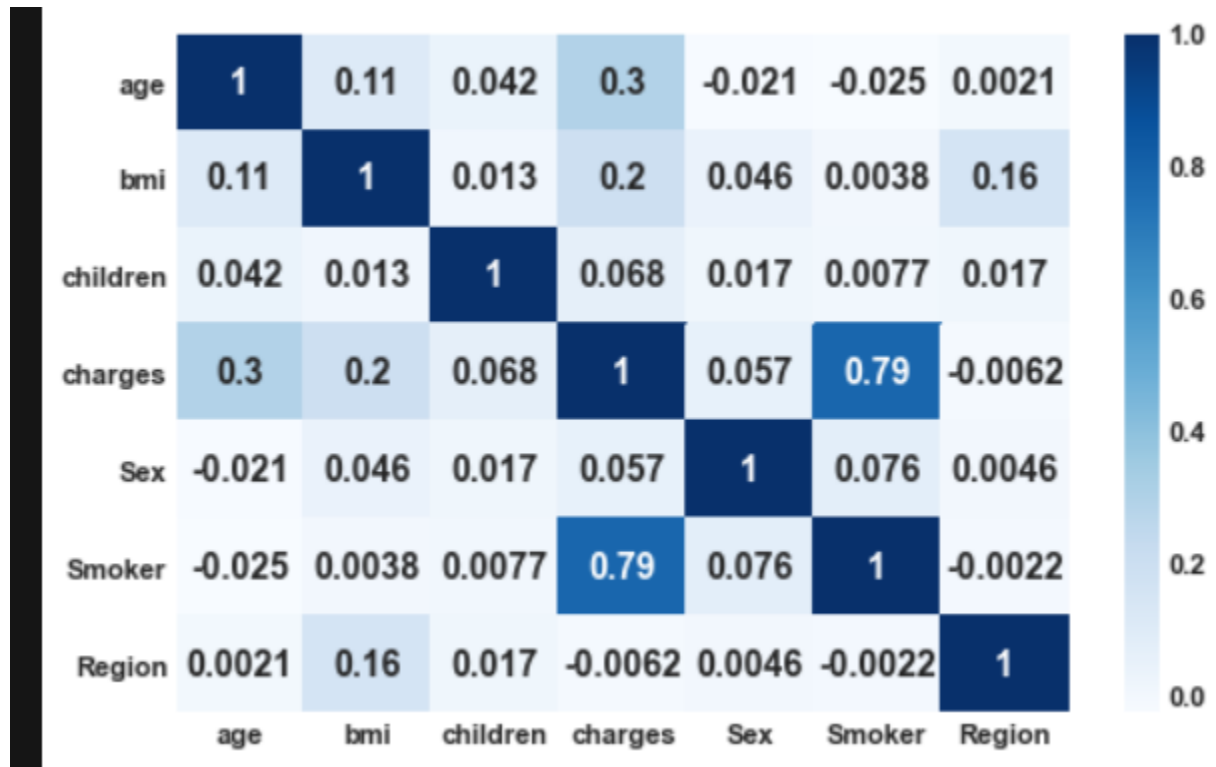
Null values:

```
age         0
sex         0
bmi         0
children    0
smoker      0
region      0
charges     0
Sex         0
Smoker      0
Region      0
```
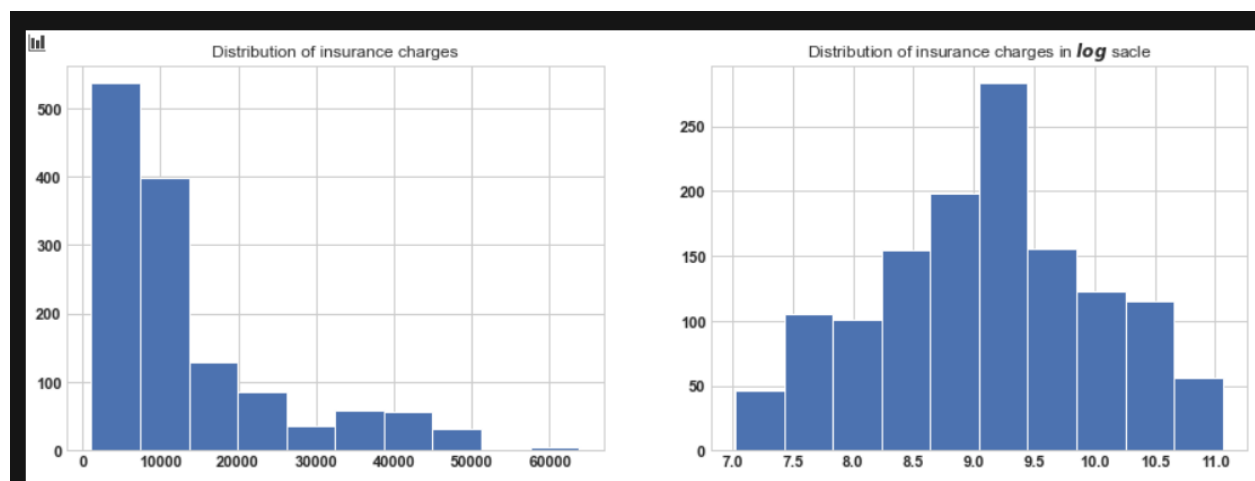
## 2)Correlation between variables

Correlations between the variables were calculated and a heatmap was plotted.

## 3)Distribution of dependent variable and its log value  was plotted

The plots show that the distribution is right skewed.



## 4)Conversion of categorical values into numerical

Categorical variables were converted into numerical values using Label encoder.

5)Data was splitted into train and test dataset with test size 0.3

6)Built the regression model using the given equation without any in-built library except numpy

Regression model was built using the equation $\theta = (X^T X)^{-1} X^T y$ using numpy functionalities.
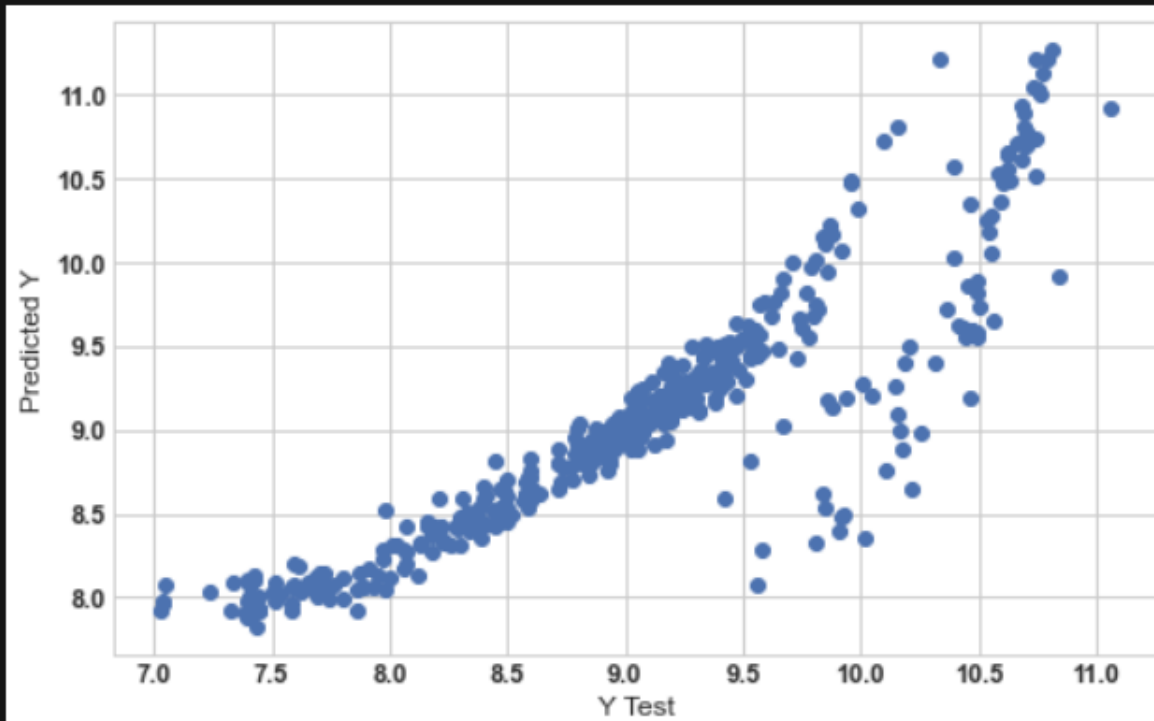
7)Trained the regression model using sklearn library

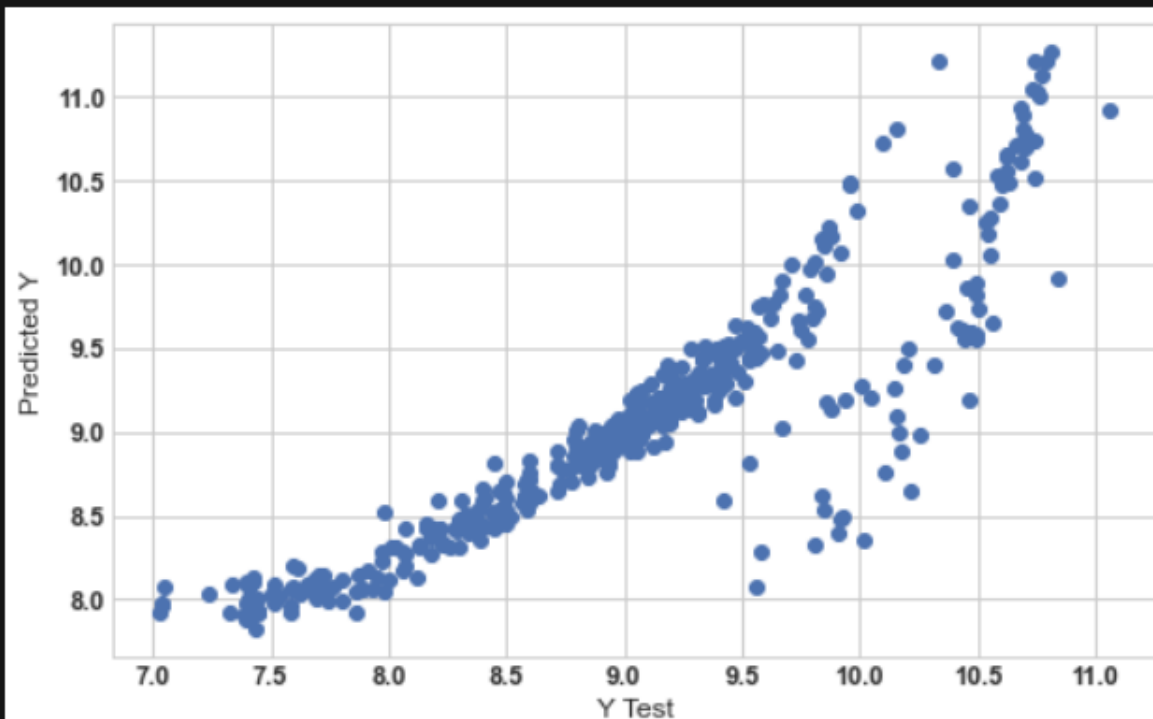8)Comparison of  parameters (coefficients) of both the models:

|  | age | bmi | children | sex | smoker | region |
|---|---|---|---|---|---|---|
| Manual | 0.03325052 | 0.01407749 | 0.10057153 | -0.04105018 | 1.52389362 | -0.05391283 |
| Using sklearn | 0.03325052 | 0.01407749 | 0.10057153 | -0.04105018 | 1.52389362 | -0.05391283 |

9)Prediction were made using both the models:

Score for manual model: 0.7892774635828786

Score for sklearn model: 0.7892774635828786

10)MSE was calculated using for both the models using manually made function and inbuilt function:

MSE using manual function :

```
The Mean Square Error(MSE) or J(theta) is:  0.17484377359606448
```

MSE using inbuilt function from sklearn:

```
The Mean Square Error(MSE) or J(theta) is:  0.17484377359606262
```

11)Plotted the relationship between actual and predicted value: